

An Approximation of the CP-Rank of a Partially Sampled Tensor

Morteza Ashraphijuo
Dep. of Electrical Eng.
Columbia University
New York, NY

Email: ashraphijuo@ee.columbia.edu

Xiaodong Wang
Dep. of Electrical Eng.
Columbia University
New York, NY

Email: wangx@ee.columbia.edu

Vaneet Aggarwal
School of Industrial Eng.
Purdue University
West Lafayette, IN

Email: vaneet@purdue.edu

Abstract—We exploit the recent algebraic geometry analyses that study the fundamental conditions on the locations of the sampled entries for finite completability of low-rank sampled tensor to treat the problem of CP-rank approximation for a partially sampled tensor. Particularly, the goal is to approximate the unknown rank based on the locations of the sampled entries, i.e., the sampling pattern, and the rank of an arbitrary given completion. First we provide an upper bound on the unknown CP-rank with probability one assuming that the sampling pattern satisfies the proposed combinatorial properties. However, the proposed combinatorial properties may be hard to verify. Hence, we also provide probabilistic versions of such bounds that hold with high probability assuming that the sampling probability is above a threshold, i.e., we provide the sampling probability that results the sampling pattern satisfies the proposed combinatorial properties with high probability. In addition, these upper bounds can be exactly equal to the unknown CP-rank given the lowest-rank completion. To illustrate how tight our proposed upper bounds are, we have provided some numerical results for the case of two-way tensor, i.e., matrix, in which we applied the nuclear norm minimization to find a low-rank completion of the sampled data and observe that the proposed upper bound is almost equal to the true unknown rank.

Index Terms—Low-rank data completion, rank estimation, tensor, matrix, manifold, CP rank.

I. INTRODUCTION

The low-rank data completion problem is concerned with completing a matrix or tensor given a subset of its entries and its rank. Various applications can be found in many fields including image and signal processing [1, 2], data mining [3], network coding [4], compressed sensing [5–7], reconstructing the visual data [8], bioinformatics [9], fingerprinting [10], systems biology [11], etc. There is an extensive literature on developing various optimization methods to treat this problem including minimizing a convex relaxation of rank [7, 12–17], non-convex approaches [18], and alternating minimization [19, 20], etc. More recently, deterministic conditions on the sampling patterns have been studied for subspace clustering in [21–24]. Moreover, the fundamental conditions on the sampling pattern that lead to different numbers of completion (unique, finite, or infinite) for different data

structures given the corresponding rank constraints have been investigated in [25–31].

However, in many practical low-rank data completion problems, the rank may not be known *a priori*. In this paper, we investigate this problem and we aim to approximate the rank based on the given entries, where it is assumed that the original data is generically chosen from the manifold corresponding to the unknown rank. The only existing work that treats this problem for matrix based on the sampling pattern is [32], which requires some strong assumptions including the existence of a completion whose rank r is a lower bound on the unknown true rank r^* , i.e., $r^* \geq r$. We start by investigating the problem for matrix or two-way tensor to provide a new analysis that does not require such assumption and also we can extend our novel approach to treat the same problem for a d -way tensor to determine its CP rank. We also obtain such bound that holds with high probability based on the sampling probability. Moreover, for the case of matrix or $d = 2$, we provide some numerical results to show how tight our probabilistic bounds on the rank are (in terms of the sampling probability). In particular, we used nuclear norm minimization to find a completion and demonstrate our proposed method in obtaining a tight bound on the unknown rank.

We take advantage of the geometric analysis on the manifold of the corresponding data which leads to the fundamental conditions on the sampling pattern (independent of the value of entries) [25, 27] such that given an arbitrary low-rank completion we can provide a tight upper bound on the rank. To illustrate how such approximation is even possible consider the following example. Assume that an $n_1 \times n_2$ rank-2 matrix is chosen generically from the corresponding manifold. Hence, any 2×2 submatrix of this matrix is full-rank with probability one (due to the genericity assumption). Moreover, note that any 3×3 submatrix of this matrix is not full-rank. As a result, by observing the sampled entries we can find some bounds on the rank. Using the analysis in [25–29, 33] on finite completability of the sampled data (finite number of completions) for different data models, we characterize both deterministic and probabilistic bounds on the unknown rank.

II. NOTATIONS AND PROBLEM STATEMENT

A. Two-Way Tensor Scenario

Assume that the sampled matrix \mathbf{U} is chosen generically from the manifold of the $n_1 \times n_2$ matrices of rank r^* , where r^* is unknown. The matrix $\mathbf{V} \in \mathbb{R}^{n_1 \times r^*}$ is called a basis for \mathbf{U} if each column of \mathbf{U} can be written as a linear combination of the columns of \mathbf{V} . Denote Ω as the binary sampling pattern matrix that is of the same size as \mathbf{U} and $\Omega(\vec{x}) = 1$ if $\mathbf{U}(\vec{x})$ is observed and $\Omega(\vec{x}) = 0$ otherwise, where $\vec{x} = (x_1, x_2)$ represents the entry corresponding to row number x_1 and column number x_2 . For each submatrix \mathbf{U}' of the matrix \mathbf{U} , define $N_\Omega(\mathbf{U}')$ as the number of observed entries in \mathbf{U}' according to the sampling pattern Ω . Moreover, define \mathbf{U}_Ω as the matrix obtained from sampling \mathbf{U} according to Ω , i.e.,

$$\mathbf{U}_\Omega(\vec{x}) = \begin{cases} \mathbf{U}(\vec{x}) & \text{if } \Omega(\vec{x}) = 1, \\ 0 & \text{if } \Omega(\vec{x}) = 0. \end{cases} \quad (1)$$

B. d -Way Tensor Scenario

Assume that a d -way tensor $\mathcal{U} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is sampled. Denote Ω as the binary sampling pattern tensor that is of the same size as \mathcal{U} and $\Omega(\vec{x}) = 1$ if $\mathcal{U}(\vec{x})$ is observed and $\Omega(\vec{x}) = 0$ otherwise, where $\mathcal{U}(\vec{x})$ represents an entry of tensor \mathcal{U} with coordinate $\vec{x} = (x_1, \dots, x_d)$. Moreover, define \mathcal{U}_Ω as the tensor obtained from sampling \mathcal{U} according to Ω , i.e.,

$$\mathcal{U}_\Omega(\vec{x}) = \begin{cases} \mathcal{U}(\vec{x}) & \text{if } \Omega(\vec{x}) = 1, \\ 0 & \text{if } \Omega(\vec{x}) = 0. \end{cases} \quad (2)$$

For each subtensor \mathcal{U}' of the tensor \mathcal{U} , define $N_\Omega(\mathcal{U}')$ as the number of observed entries in \mathcal{U}' according to the sampling pattern Ω .

The CP rank of a tensor \mathcal{U} , $\text{rank}_{\text{CP}}(\mathcal{U}) = r$, is defined as the minimum number r such that there exist $\mathbf{a}_i^l \in \mathbb{R}^{n_i}$ for $1 \leq i \leq d$ and $1 \leq l \leq r$, such that

$$\mathcal{U} = \sum_{l=1}^r \mathbf{a}_1^l \otimes \mathbf{a}_2^l \otimes \dots \otimes \mathbf{a}_d^l, \quad (3)$$

or equivalently,

$$\mathcal{U}(x_1, x_2, \dots, x_d) = \sum_{l=1}^r \mathbf{a}_1^l(x_1) \mathbf{a}_2^l(x_2) \dots \mathbf{a}_d^l(x_d), \quad (4)$$

where \otimes denotes the tensor product (outer product) and $\mathbf{a}_i^l(x_i)$ denotes the x_i -th entry of vector \mathbf{a}_i^l . Note that $\mathbf{a}_1^l \otimes \mathbf{a}_2^l \otimes \dots \otimes \mathbf{a}_d^l \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is a rank-1 tensor, $l = 1, 2, \dots, r$.

C. Problem Statement

In this paper, we assume that there exists a full rank completion of the sampled data (i.e., the data is not over-sampled). We are interested in obtaining the upper bound on the unknown rank r^* deterministically based on the sampling pattern Ω or Ω and the rank of a given completion. Also, we aim to provide such bound that holds with high probability based only on the sampling probability of the entries and the rank of a given completion. Moreover, we

provide both deterministic and probabilistic conditions such that the unknown rank can be exactly determined.

III. MAIN RESULTS

Note that for the case of two-way tensor or matrix, the CP rank is simply equivalent to the well-known matrix rank. Hence, for the simplicity we start by proposing our approach for the scenario of $d = 2$ in Section III-A and then generalize it to a general d -way tensor in Section III-B.

A. Matrix Analysis

Previously, this problem has been treated in [32], where strong assumptions including the existence of a completion with rank $r \leq r^*$ have been used. In this section, we provide an analysis that does not require such assumption and moreover our analysis can be extended to tensors as we show in Section III-B. Furthermore, we show the tightness of our theoretical bounds via numerical examples.

1) Deterministic Rank Analysis:

The following assumption will be used frequently in this subsection.

Assumption A_r : Each column of the sampled matrix includes at least r sampled entries.

Consider an arbitrary column of the sampled matrix $\mathbf{U}(:, i)$, where $i \in \{1, \dots, n_2\}$. Let $l_i = N_\Omega(\mathbf{U}(:, i))$ denote the number of observed entries in the i -th column of \mathbf{U} . Assumption A_r results that $l_i \geq r$.

We construct a binary valued matrix called **constraint matrix** $\check{\Omega}_r$ based on Ω and a given number r . Specifically, we construct $l_i - r$ columns with binary entries based on the locations of the observed entries in $\mathbf{U}(:, i)$ such that each column has exactly $r + 1$ entries equal to one. Assume that x_1, \dots, x_{l_i} are the row indices of all observed entries in this column. Let Ω_r^i be the corresponding $n_1 \times (l_i - r)$ matrix to this column which is defined as the following: for any $j \in \{1, \dots, l_i - r\}$, the j -th column has the value 1 in rows $\{x_1, \dots, x_r, x_{r+j}\}$ and zeros elsewhere. Define the binary constraint matrix as $\check{\Omega}_r = [\Omega_r^1 | \Omega_r^2 \dots | \Omega_r^{n_2}] \in \mathbb{R}^{n_1 \times K_r}$ [25], where $K_r = N_\Omega(\mathbf{U}) - n_2 r$.

Assumption B_r : There exists a submatrix¹ $\check{\Omega}_r' \in \mathbb{R}^{n_1 \times K}$ of $\check{\Omega}_r$ such that $K = n_1 r - r^2$ and for any $K' \in \{1, 2, \dots, K\}$ and any submatrix $\check{\Omega}_r'' \in \mathbb{R}^{n_1 \times K'}$ of $\check{\Omega}_r'$ we have

$$r f(\check{\Omega}_r'') - r^2 \geq K', \quad (5)$$

where $f(\check{\Omega}_r'')$ denotes the number of nonzero rows of $\check{\Omega}_r''$.

Note that exhaustive enumeration is needed in order to check whether or not Assumption B_r holds. Hence, the deterministic analysis cannot be used in practice for large-scale data. However, it serves as the basis of the subsequent probabilistic analysis that will lead to a simple lower bound on the sampling probability such that Assumption B_r holds with high probability, which is of practical value.

¹Specified by a subset of rows and a subset of columns (not necessarily consecutive).

In the following, we restate Theorem 1 in [25] which will be used later.

Lemma 1. *For almost every \mathbf{U} , there are finitely many completions of the sampled matrix if and only if Assumptions A_{r^*} and B_{r^*} hold.*

Recall that the true rank r^* is assumed unknown.

Definition 1. *Let \mathcal{S}_Ω denote the set of all natural numbers r such that both Assumptions A_r and B_r hold.*

Lemma 2. *There exists a number r_Ω such that $\mathcal{S}_\Omega = \{1, 2, \dots, r_\Omega\}$.*

Proof. Assume that $1 < r \leq \min\{n_1, n_2\}$ and $r \in \mathcal{S}_\Omega$. It suffices to show $r - 1 \in \mathcal{S}_\Omega$. By contradiction, assume that $r - 1 \notin \mathcal{S}_\Omega$. Therefore, according to Lemma 1, there exist infinitely many completions of \mathbf{U} of rank $r - 1$ and there exist at most finitely many completions of \mathbf{U} of rank r .

Consider a rank $r - 1$ completion \mathbf{U}_{r-1} . Note that changing one single entry (a non-observed entry) of \mathbf{U}_{r-1} , for example $\mathbf{U}_{r-1}(1, 1) = x$, to a random number in $y \in \mathbb{R}$ changes the rank of this matrix by at most 1 and basically since we are changing to a random number, it can be easily seen that the rank does not decrease with probability one. Hence, the rank of the modified matrix \mathbf{U}'_{r-1} would be either $r - 1$ or r . Assume that the rank has been increased to r . Then, we show there exist infinitely many completions of rank r , which contradicts the assumption. In fact, for any value of $\mathbf{U}_{r-1}(1, 1)$ except x , this matrix would be a rank r completion. To observe this more clearly, consider the $r \times r$ submatrix of \mathbf{U}'_{r-1} whose determinant is not zero due to changing the value of $\mathbf{U}_{r-1}(1, 1)$. It is easily observed that this submatrix includes $\mathbf{U}'_{r-1}(1, 1)$ and let assume it is $\mathbf{U}'_{r-1}(1 : r, 1 : r)$, and therefore the determinant of $\mathbf{U}'_{r-1}(2 : r, 2 : r)$ is a nonzero number (otherwise the rank would not increase by changing the value of $\mathbf{U}_{r-1}(1, 1)$). Hence, it is easy to see that for any value of $\mathbf{U}_{r-1}(1, 1)$ except x , \mathbf{U}'_{r-1} would be a rank r completion, and therefore there exist infinitely many completions of rank r and proof is complete in this scenario.

Now, assume that changing any of the non-observed entries does not increase the rank of \mathbf{U}_{r-1} . Then, this contradicts the assumption that there exists a full rank completion of the sampled data since there does not exist any completion of rank higher than $r - 1$. \square

The following theorem provides a relationship between the unknown rank r^* and r_Ω .

Theorem 1. *With probability one, exactly one of the following statements holds*

(i) $r^* \in \mathcal{S}_\Omega = \{1, 2, \dots, r_\Omega\}$;

(ii) *For any arbitrary completion of the sampled matrix \mathbf{U} of rank r , we have $r \notin \mathcal{S}_\Omega$.*

Proof. Suppose that there does not exist a completion of the sampled matrix \mathbf{U} of rank r such that $r \in \mathcal{S}_\Omega$. Therefore, it

is easily verified that statement (ii) holds and statement (i) does not hold. On the other hand, assume that there exists a completion of the sampled matrix \mathbf{U} of rank r , where $r \in \mathcal{S}_\Omega$. Hence, statement (ii) does not hold and to complete the proof it suffices to show that with probability one, statement (i) holds.

Observe that $r_\Omega \in \mathcal{S}_\Omega$, and therefore Assumption A_{r_Ω} holds. Hence, each column of \mathbf{U} includes at least $r_\Omega + 1$ observed entries. On the other hand, the existence of a completion of the sampled matrix \mathbf{U} of rank $r \in \mathcal{S}_\Omega$ results in the existence of a basis $\mathbf{X} \in \mathbb{R}^{n_1 \times r}$ such that each column of \mathbf{U} is a linear combination of the columns of \mathbf{X} , and thus there exists $\mathbf{Y} \in \mathbb{R}^{r \times n_2}$ such that $\mathbf{U}_\Omega = (\mathbf{X}\mathbf{Y})_\Omega$. Hence, given \mathbf{X} , each observed entry $\mathbf{U}(i, j)$ results in a degree-1 polynomial in terms of the entries of \mathbf{Y} as the following

$$\mathbf{U}(i, j) = \sum_{l=1}^r \mathbf{X}(i, l)\mathbf{Y}(l, j). \quad (6)$$

Consider the first column of \mathbf{U} and recall that it includes at least $r_\Omega + 1 \geq r + 1$ observed entries. The genericity of the coefficients of the above-mentioned polynomials results that using r of the observed entries the first column of \mathbf{Y} can be determined uniquely. This is because there exists a unique solution for a system of r linear equations in r variables that are linearly independent. Then, there exists at least one more observed entry besides these r observed entries in the first column of \mathbf{U} and it can be written as a linear combination of the r observed entries that have been used to obtain the first column of \mathbf{Y} . Let $\mathbf{U}(i_1, 1), \dots, \mathbf{U}(i_r, 1)$ denote the r observed entries that have been used to obtain the first column of \mathbf{Y} and $\mathbf{U}(i_{r+1}, 1)$ denote the other observed entry. Hence, the existence of a completion of the sampled matrix \mathbf{U} of rank $r \in \mathcal{S}_\Omega$ results in an equation as the following

$$\mathbf{U}(i_{r+1}, 1) = \sum_{l=1}^r t_l \mathbf{U}(i_l, 1), \quad (7)$$

where t_l 's are constant scalars, $l = 1, \dots, r$. Assume that $r^* \notin \mathcal{S}_\Omega$, i.e., statement (i) does not hold. Then, note that $r^* \geq r + 1$ and \mathbf{U} is chosen generically from the manifold of $n_1 \times n_2$ rank- r^* matrices, and therefore an equation of the form of (7) holds with probability zero. Moreover, according to Lemma 1 there exist at most finitely many completions of the sampled matrix of rank r . Therefore, there exist a completion of \mathbf{U} of rank r with probability zero, which contradicts the initial assumption that there exists a completion of the sampled matrix \mathbf{U} of rank r , where $r \in \mathcal{S}_\Omega$. \square

Corollary 1. *Consider an arbitrary number $r' \in \mathcal{S}_\Omega$. Similar to Theorem 1, it follows that with probability one, exactly one of the followings holds*

(i) $r^* \in \{1, 2, \dots, r'\}$;

(ii) *For any arbitrary completion of the sampled matrix \mathbf{U} of rank r , we have $r \notin \{1, 2, \dots, r'\}$.*

As a result of Corollary 1, we have the following.

Corollary 2. Assuming that there exists a rank- r completion of the sampled matrix \mathbf{U} such that $r \in \mathcal{S}_\Omega$, then with probability one $r^* \leq r$.

Corollary 3. Let \mathbf{U}^* denote an optimal solution to the following NP-hard optimization problem

$$\begin{aligned} & \text{minimize}_{\mathbf{U}' \in \mathbb{R}^{n_1 \times n_2}} \quad \text{rank}(\mathbf{U}') \\ & \text{subject to} \quad \mathbf{U}'_\Omega = \mathbf{U}_\Omega. \end{aligned} \quad (8)$$

Also, let $\hat{\mathbf{U}}$ denote a suboptimal solution to the above optimization problem. Then, Corollary 1 results the following statements:

- (i) If $\text{rank}(\mathbf{U}^*) \in \mathcal{S}_\Omega$, then $r^* = \text{rank}(\mathbf{U}^*)$ with probability one.
- (ii) If $\text{rank}(\hat{\mathbf{U}}) \in \mathcal{S}_\Omega$, then $r^* \leq \text{rank}(\hat{\mathbf{U}})$ with probability one.

Remark 1. One challenge of applying Corollary 3 or any of the other obtained deterministic results is the computation of \mathcal{S}_Ω , which involves exhaustive enumeration to check Assumption B_r. Next, for each number r , we provide a lower bound on the sampling probability in terms of r that ensures $r \in \mathcal{S}_\Omega$ with high probability. Consequently, we do not need to compute \mathcal{S}_Ω but instead we can certify the above results with high probability.

2) Probabilistic Rank Analysis:

The following lemma is a re-statement of Theorem 3 in [25], which is the probabilistic version of Lemma 1.

Lemma 3. Suppose $r \leq \frac{n_1}{6}$ and that each column of the sampled matrix is observed in at least l entries, uniformly at random and independently across entries, where

$$l > \max \left\{ 12 \log \left(\frac{n_1}{\epsilon} \right) + 12, 2r \right\}. \quad (9)$$

Also, assume that $r(n_1 - r) \leq n_2$. Then, with probability at least $1 - \epsilon$, $r \in \mathcal{S}_\Omega$.

The following lemma is taken from [28] and will be used to derive a lower bound on the sampling probability that leads to the similar statement as Theorem 1 with high probability.

Lemma 4. Consider a vector with n entries where each entry is observed with probability p independently from the other entries. If $p > p' = \frac{k}{n} + \frac{1}{\sqrt[4]{n}}$, then with probability at least $\left(1 - \exp(-\frac{\sqrt{n}}{2})\right)$, more than k entries are observed.

The following proposition characterizes the probabilistic version of Theorem 1.

Proposition 1. Suppose $r \leq \frac{n_1}{6}$, $r(n_1 - r) \leq n_2$ and that each entry of the sampled matrix is observed uniformly at random and independently across entries with probability p , where

$$p > \frac{1}{n_1} \max \left\{ 12 \log \left(\frac{n_1}{\epsilon} \right) + 12, 2r \right\} + \frac{1}{\sqrt[4]{n_1}}. \quad (10)$$

Then, with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$, we have $r \in \mathcal{S}_\Omega$.

Proof. Consider an arbitrary column of \mathbf{U} and note that resulting from Lemma 4 the number of observed entries at this column of \mathbf{U} is greater than $\max \left\{ 12 \log \left(\frac{n_1}{\epsilon} \right) + 12, 2r \right\}$ with probability at least $\left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)$. Therefore, the number of sampled entries at each column satisfies

$$l > \max \left\{ 12 \log \left(\frac{n_1}{\epsilon} \right) + 12, 2r \right\}, \quad (11)$$

with probability at least $\left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$. Thus, resulting from Lemma 3 with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$, we have $r \in \mathcal{S}_\Omega$. \square

Finally, we have the following probabilistic version of Corollary 3.

Corollary 4. Assume that $\text{rank}(\mathbf{U}^*) \leq \frac{n_1}{6}$ and $\text{rank}(\mathbf{U}^*)(n_1 - \text{rank}(\mathbf{U}^*)) \leq n_2$ and (10) holds for $r = \text{rank}(\mathbf{U}^*)$, where \mathbf{U}^* denotes an optimal solution to the optimization problem (8). Then, according to Proposition 1 and Corollary 3, with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$, $r^* = \text{rank}(\mathbf{U}^*)$. Similarly, assume that $\text{rank}(\hat{\mathbf{U}}) \leq \frac{n_1}{6}$ and $\text{rank}(\hat{\mathbf{U}})(n_1 - \text{rank}(\hat{\mathbf{U}})) \leq n_2$ and (10) holds for $r = \text{rank}(\hat{\mathbf{U}})$, where $\hat{\mathbf{U}}$ denotes a suboptimal solution to the optimization problem (8). Then, with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$, $r^* \leq \text{rank}(\hat{\mathbf{U}})$.

3) Numerical Results:

In Fig. 1 and Fig. 2, the x-axis represents the sampling probability, and the y-axis denotes the value of r . The color scale represents the lower bound on the probability of event $r \in \mathcal{S}_\Omega$. For example, as we can observe in Fig. 1, for any $r \in \{1, \dots, 44\}$ we have $r \in \mathcal{S}_\Omega$ with probability at least 0.6 (approximately based on the color scale since the corresponding points are orange) given that $p = 0.54$.

We consider the sampled matrix $\mathbf{U} \in \mathbb{R}^{300 \times 15000}$ and $\mathbf{U} \in \mathbb{R}^{1200 \times 240000}$ in Fig. 1 and Fig. 2, respectively. In particular, for fixed values of sampling probability p and r , we first find a “small” ϵ that (10) holds by trial-and-error. Then, according to Proposition 1, we conclude that with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n_1}}{2})\right)^{n_2}$, $r \in \mathcal{S}_\Omega$.

The purpose of Fig. 3 is to show how tight our proposed upper bounds on rank can be. Here, we first generate an $n_1 \times n_2$ random matrix of a given rank r by multiplying a random (entries are drawn according to a uniform distribution on real numbers within an interval) $n_1 \times r$ matrix and $r \times n_2$ matrix. Then, each entry of the randomly generated matrix is sampled uniformly at random and independently across entries with some sampling probability p . Afterwards, we apply the nuclear norm minimization method proposed in [36] for matrix completion, where the non-convex objective function in (8) is relaxed by using nuclear norm, which is

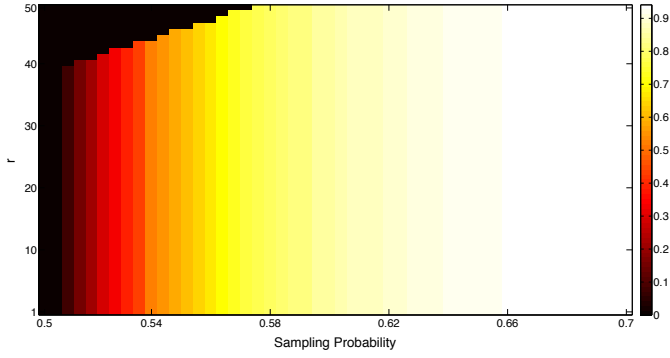


Fig. 1: Probability of $r \in \mathcal{S}_\Omega$ as a function of sampling probability for $\mathbf{U} \in \mathbb{R}^{300 \times 15000}$.

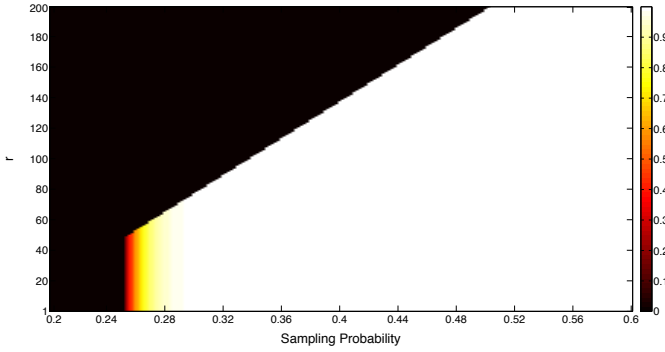


Fig. 2: Probability of $r \in \mathcal{S}_\Omega$ as a function of sampling probability for $\mathbf{U} \in \mathbb{R}^{1200 \times 240000}$.

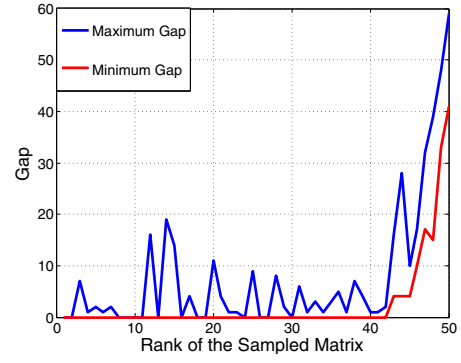
the convex hull of the rank function, as follows

$$\begin{aligned} & \text{minimize}_{\mathbf{U}' \in \mathbb{R}^{n_1 \times n_2}} \quad \|\mathbf{U}'\|_* \\ & \text{subject to} \quad \mathbf{U}'_\Omega = \mathbf{U}_\Omega, \end{aligned} \quad (12)$$

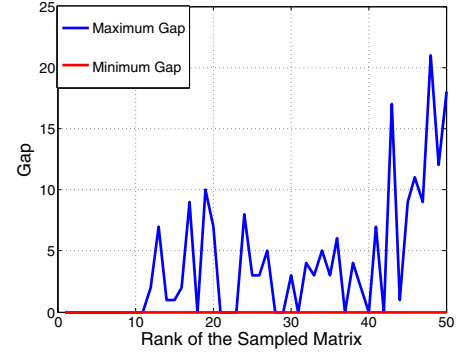
where $\|\mathbf{U}'\|_*$ denotes the nuclear norm of \mathbf{U}' . Let $\hat{\mathbf{U}}^*$ denote an optimal solution to (12) and recall that \mathbf{U}^* denotes an optimal solution to (8). Since (12) is a convex relaxation to (8), we conclude that $\hat{\mathbf{U}}^*$ is a suboptimal solution to (8), and therefore $\text{rank}(\mathbf{U}^*) \leq \text{rank}(\hat{\mathbf{U}}^*)$. We used the Matlab program found online [37] to solve (12).

As an example, we generate a random matrix $\mathbf{U} \in \mathbb{R}^{300 \times 15000}$ (the same size as the matrix in Fig. 1) of rank r as described above for $r \in \{1, \dots, 50\}$ and some values of the sampling probability p . Then, we obtain the rank of the completion given by (12) and denote it by r' . Due to the randomness of the sampled matrix, we repeat this procedure 5 times. We calculate the “gap” $r' - r$ in each of these 5 runs and denote the maximum and minimum among these 5 numbers by d_{\max} and d_{\min} , respectively. Hence, d_{\max} and d_{\min} represent the loosest (worst) and tightest (best) gaps between the rank obtained by (12) and rank of the original sampled matrix over 5 runs, respectively. In Fig. 3, the maximum and minimum gaps are plotted as a function of rank of the matrix, for different sampling probabilities.

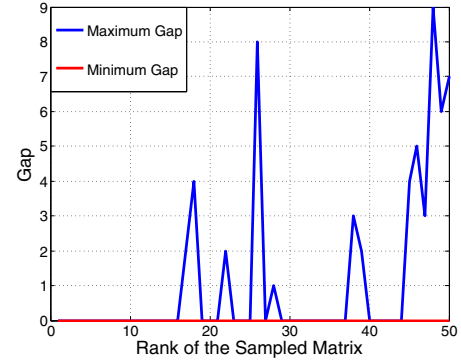
We have the following observations.



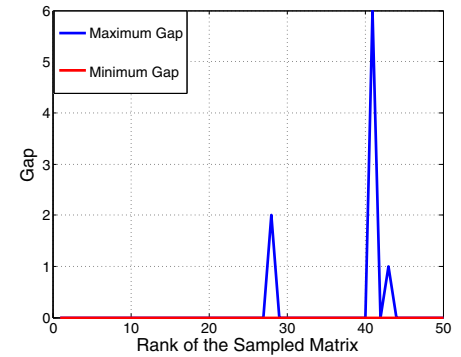
(a) $p = 0.46$.



(b) $p = 0.50$.



(c) $p = 0.54$.



(d) $p = 0.58$.

Fig. 3: The gaps between the rank of the obtained matrix via (12) and that of the original sampled matrix.

- According to Fig. 1, for $p = 0.54$ and $p = 0.58$ we can ensure that the rank of any completion is an upper bound on the rank of the sampled matrix or r^* with probability at least 0.6 and 0.8, respectively.
- As we can observe in Fig. 3(a)-(d), the defined gap is always a nonnegative number, which is consistent with previous observation that for $p = 0.54$ and $p = 0.58$ we can certify that with high probability (≥ 0.6) the rank of any completion is an upper bound on the rank of the sampled matrix or r^* .
- For $p = 0.54$ and $p = 0.58$ that we have theoretical results (as mentioned in the first observation) the gap obtained by (12) is very close to zero. This phenomenon (that we do not have a rigorous justification for) shows that as soon as we can certify our proposed theoretical results (i.e., as soon as the rank of a completion provides an upper bound on the rank of the sampled matrix or r^*) by increasing the sampling probability, the upper bound found through (12) becomes very tight; in some cases this bound is exactly equal to r^* (red curves) and in some cases this bound is almost equal to r^* (blue curves). However, these gaps are not small (specially blue curves) for $p = 0.46$ and $p = 0.50$ and note that according to Fig. 1, for these values of p we cannot guarantee the bounds on the value of rank hold with high probability.

B. Tensor Analysis

In this subsection, we assume that the sampled tensor $\mathcal{U} \in \mathbb{R}^{n_1 \times \dots \times n_d}$ is chosen generically from the manifold of tensors of rank $r^* = \text{rank}_{\text{CP}}(\mathcal{U})$, where r^* is unknown.

Assumption \mathcal{A}_r : Each row of the d -th matricization of the sampled tensor, i.e., $\mathbf{U}_{(d)}$ includes at least r observed entries.

We construct a binary valued tensor called **constraint tensor** $\check{\Omega}_r$ based on Ω and a given number r . Consider any subtensor $\mathcal{Y} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times 1}$ of the tensor \mathcal{U} . The sampled tensor \mathcal{U} includes n_d subtensors that belong to $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times 1}$ and let \mathcal{Y}_i for $1 \leq i \leq n_d$ denote these n_d subtensors. Define a binary valued tensor $\check{\mathcal{Y}}_i \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times k_i}$, where $k_i = N_\Omega(\mathcal{Y}_i) - r$ and its entries are described as the following. We can look at $\check{\mathcal{Y}}_i$ as k_i tensors each belongs to $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times 1}$. For each of the mentioned k_i tensors in $\check{\mathcal{Y}}_i$ we set the entries corresponding to r of the observed entries equal to 1. For each of the other k_i observed entries, we pick one of the k_i tensors of $\check{\mathcal{Y}}_i$ and set its corresponding entry (the same location as that specific observed entry) equal to 1 and set the rest of the entries equal to 0. In the case that $k_i = 0$ we simply ignore $\check{\mathcal{Y}}_i$, i.e., $\check{\mathcal{Y}}_i = \emptyset$.

By putting together all n_d tensors in dimension d , we construct a binary valued tensor $\check{\Omega}_r \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times K}$, where $K = \sum_{i=1}^{n_d} k_i = N_\Omega(\mathcal{U}) - rn_d$ and call it the **constraint tensor**. Observe that each subtensor of $\check{\Omega}_r$ which belongs to $\mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times 1}$ includes exactly $r+1$ nonzero entries. In [27], an example is given on the construction of $\check{\Omega}_r$.

Assumption \mathcal{B}_r : $\check{\Omega}_r$ consists a subtensor $\check{\Omega}'_r \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times K}$ such that $K = r(\sum_{i=1}^{d-1} n_i) - r^2 - r(d-2)$ and for any $K' \in \{1, 2, \dots, K\}$ and any subtensor $\check{\Omega}''_r \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_{d-1} \times K'}$ of $\check{\Omega}'_r$ we have

$$r \left(\sum_{i=1}^{d-1} f_i(\check{\Omega}''_r) \right) - r \min \left\{ \max \left\{ f_1(\check{\Omega}''_r), \dots, f_{d-1}(\check{\Omega}''_r) \right\}, r \right\} - (d-2) \geq K',$$

where $f_i(\check{\Omega}''_r)$ denotes the number of nonzero rows of the i -th matricization of $\check{\Omega}''_r$.

The following lemma is a re-statement of Theorem 1 in [27].

Lemma 5. *For almost every \mathcal{U} , there are only finitely many rank- r^* completions of the sampled tensor if and only if Assumptions \mathcal{A}_{r^*} and \mathcal{B}_{r^*} hold.*

Definition 2. *Let \mathcal{S}_Ω denote the set of all natural numbers r such that both Assumptions \mathcal{A}_r and \mathcal{B}_r hold.*

Lemma 6. *There exists a number r_Ω such that $\mathcal{S}_\Omega = \{1, 2, \dots, r_\Omega\}$.*

Proof. The proof is similar to the proof of Lemma 2 since the dimension of the manifold of CP rank- r tensors is $r(\sum_{i=1}^d n_i) - r^2 - r(d-1)$, which is an increasing function in r . \square

The following theorem gives an upper bound on the unknown rank r^* .

Theorem 2. *For almost every \mathcal{U} , with probability one, exactly one of the following statements holds*

- $r^* \in \mathcal{S}_\Omega = \{1, 2, \dots, r_\Omega\}$;
- For any arbitrary completion of the sampled tensor \mathcal{U} of rank r , we have $r \notin \mathcal{S}_\Omega$.

Proof. Similar to the proof of Theorem 1, it suffices to show that the assumption $r^* \notin \mathcal{S}_\Omega$ results that there exists a completion of \mathcal{U} of CP rank r , where $r \in \mathcal{S}_\Omega$, with probability zero. Define $\mathcal{V} = (\mathcal{V}_1, \dots, \mathcal{V}_r)$ as the basis of the rank- r CP decomposition of \mathcal{U} as in (3), where $\mathcal{V}_l = \mathbf{a}_1^l \otimes \mathbf{a}_2^l \otimes \dots \otimes \mathbf{a}_{d-1}^l \in \mathbb{R}^{n_1 \times \dots \times n_{d-1}}$ is a rank-1 tensor and \mathbf{a}_l^l is defined in (3) for $1 \leq l \leq r$ and $1 \leq i \leq d$. Define $\mathcal{Y} = (\mathbf{a}_d^1, \dots, \mathbf{a}_d^r)$ and $\mathcal{V} \otimes_d \mathcal{Y} = \sum_{l=1}^r \mathcal{V}_l \otimes \mathbf{a}_d^l$. Observe that $\mathcal{U} = \sum_{l=1}^r \mathcal{V}_l \otimes \mathbf{a}_d^l = \mathcal{V} \otimes_d \mathcal{Y}$.

Observe that each row of $\mathbf{U}_{(d)}$ includes at least $r_\Omega + 1$ observed entries since Assumption \mathcal{A}_{r_Ω} holds. Moreover, the existence of a completion of the sampled tensor \mathcal{U} of rank $r \in \mathcal{S}_\Omega$ results in the existence of a basis $\mathcal{V} = (\mathcal{V}_1, \dots, \mathcal{V}_r)$ such that there exists $\mathcal{Y} = (\mathbf{a}_d^1, \dots, \mathbf{a}_d^r)$ and $\mathcal{U}_\Omega = (\mathcal{V} \otimes_d \mathcal{Y})_\Omega$. As a result, given \mathcal{V} , each observed entry of \mathcal{U} results in a degree-1 polynomial in terms of the entries of \mathcal{Y} as

$$\mathcal{U}(\vec{x}) = \sum_{l=1}^r \mathcal{V}_l(x_1, \dots, x_{d-1}) \mathbf{a}_d^l(x_d). \quad (13)$$

Note that $r_\Omega \geq r$ and each row of $\mathbf{U}_{(d)}$ includes at least $r_\Omega + 1 \geq r + 1$ observed entries. Consider $r + 1$ of the

observed entries of the first row of $\mathbf{U}_{(d)}$ and we denote them by $\mathcal{U}(\vec{x}_1), \dots, \mathcal{U}(\vec{x}_{r+1})$, where the last component of the vector \vec{x}_i is equal to one, $1 \leq i \leq r+1$. Similar to the proof of Theorem 1, genericity of \mathcal{U} results in

$$\mathcal{U}(\vec{x}_{r+1}) = \sum_{l=1}^r t_l \mathcal{U}(\vec{x}_l), \quad (14)$$

where t_l 's are constant scalars, $l = 1, \dots, r$. On the other hand, according to Lemma 5 there exist at most finitely many completions of the sampled tensor of rank r . Therefore, there exist a completion of \mathbf{U} of rank r with probability zero. Moreover, an equation of the form of (14) holds with probability zero as $r^* \geq r+1$ and \mathcal{U} is chosen generically from the manifold of tensors of rank- r^* . Therefore, there exists a completion of rank r with probability zero. \square

Corollary 5. Consider an arbitrary number $r' \in \mathcal{S}_\Omega$. Similar to Theorem 2, it follows that with probability one, exactly one of the followings holds

(i) $r^* \in \{1, 2, \dots, r'\}$;

(ii) For any arbitrary completion of the sampled tensor \mathcal{U} of rank r , we have $r \notin \{1, 2, \dots, r'\}$.

Corollary 6. Assuming that there exists a CP rank- r completion of the sampled tensor \mathcal{U} such that $r \in \mathcal{S}_\Omega$, we conclude that with probability one $r^* \leq r$.

Corollary 7. Let \mathcal{U}^* denote an optimal solution to the following NP-hard optimization problem

$$\begin{aligned} & \text{minimize}_{\mathcal{U}' \in \mathbb{R}^{n_1 \times \dots \times n_d}} \text{rank}_{\text{CP}}(\mathcal{U}') \\ & \text{subject to} \quad \mathcal{U}'_\Omega = \mathcal{U}_\Omega. \end{aligned} \quad (15)$$

Assume that $\text{rank}_{\text{CP}}(\mathcal{U}^*) \in \mathcal{S}_\Omega$. Then, Corollary 6 results that $r^* = \text{rank}_{\text{CP}}(\mathcal{U}^*)$ with probability one.

The following lemma is Lemma 15 in [27], which is the probabilistic version of Lemma 5 in terms of the sampling probability.

Lemma 7. Assume that $n_1 = n_2 = \dots = n_d = n$, $d > 2$, $n > \max\{200, r(d-2)\}$ and $r \leq \frac{n}{6}$. Moreover, assume that the sampling probability satisfies

$$p > \frac{1}{n^{d-2}} \max \left\{ 27 \log \left(\frac{n}{\epsilon} \right) + 9 \log \left(\frac{2r(d-2)}{\epsilon} \right) + 18, 6r \right\} + \frac{1}{\sqrt[4]{n^{d-2}}}. \quad (16)$$

Then, with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n^{d-2}}}{2}) \right)^{n^2}$, we have $r \in \mathcal{S}_\Omega$.

The following corollary is the probabilistic version of Corollaries 6 and 7.

Corollary 8. Assuming that there exists a CP rank- r completion of the sampled tensor \mathcal{U} such that the conditions given in Lemma 7 hold, with the sampling probability satisfying (16), we conclude that with probability at least

$(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n^{d-2}}}{2}) \right)^{n^2}$ we have $r^* \leq r$. Therefore, given that (16) holds for $r = \text{rank}(\mathbf{U}^*)$ and \mathbf{U}^* denotes an optimal solution to the optimization problem (15), with probability at least $(1 - \epsilon) \left(1 - \exp(-\frac{\sqrt{n^{d-2}}}{2}) \right)^{n^2}$ we have $r^* = \text{rank}(\mathbf{U}^*)$.

IV. CONCLUSIONS

Recently, fundamental conditions on the sampling patterns have been obtained for finite completability of low-rank matrices or tensors given the corresponding ranks. In this paper, we consider the scenario where the rank is not given and we aim to approximate the unknown rank based on the location of sampled entries and some given completion. We consider a tensor and provide an upper bound on the CP-rank when an arbitrary low-CP-rank completion is given. We have characterized these bounds both deterministically, i.e., with probability one given that the sampling pattern satisfies certain combinatorial properties, and probabilistically, i.e., with high probability given that the sampling probability is above some threshold. Moreover, we showed that the obtained upper bound is exactly equal to the unknown rank if the lowest-rank completion is given. Furthermore, we have provided numerical experiments for the case of two-way tensors, where we use nuclear norm minimization to find a low-rank completion of the sampled data and we observe that in most of the cases the proposed upper bound on the rank is equal to the true rank.

REFERENCES

- [1] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski, "Phase retrieval via matrix completion," *SIAM Journal on Imaging Sciences*, vol. 6, no. 1, pp. 199–225, 2013.
- [2] H. Ji, C. Liu, Z. Shen, and Y. Xu, "Robust video denoising using low rank matrix completion," in *CVPR*, 2010, pp. 1791–1798.
- [3] L. Eldén, *Matrix Methods in Data Mining and Pattern Recognition*. Society for Industrial and Applied Mathematics, 2007, vol. 4.
- [4] N. J. Harvey, D. R. Karger, and K. Murota, "Deterministic network coding by matrix completion," in *the annual ACM-SIAM symposium on discrete algorithms*, 2005, pp. 489–498.
- [5] L.-H. Lim and P. Comon, "Multirray signal processing: Tensor decomposition meets compressed sensing," *Comptes Rendus Mecanique*, vol. 338, no. 6, pp. 311–320, 2010.
- [6] N. D. Sidiropoulos and A. Kyriklidis, "Multi-way compressed sensing for sparse low-rank tensors," *IEEE Signal Processing Letters*, vol. 19, no. 11, pp. 757–760, 2012.
- [7] S. Gandy, B. Recht, and I. Yamada, "Tensor completion and low-n-rank tensor recovery via convex optimization," *Inverse Problems*, vol. 27, no. 2, pp. 1–19, 2011.
- [8] J. Liu, P. Musialski, P. Wonka, and J. Ye, "Tensor completion for estimating missing values in visual data," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 208–220, 2013.
- [9] O. E. Ogundijo, A. Elmas, and X. Wang, "Reverse engineering gene regulatory networks from measurement with missing values," *EURASIP Journal on Bioinformatics and Systems Biology*, vol. 2017, no. 1, p. 2, 2017.
- [10] X. Y. Liu, S. Aeron, V. Aggarwal, X. Wang, and M. Y. Wu, "Adaptive sampling of RF fingerprints for fine-grained indoor localization," *IEEE Transactions on Mobile Computing*, vol. 15, no. 10, pp. 2411–2423, 2016.
- [11] O. E. Ogundijo, A. Elmas, and X. Wang, "Supplemental material for reverse engineering gene regulatory networks from measurement with missing values."

- [12] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Foundations of Computational Mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [13] E. J. Candès and T. Tao, "The power of convex relaxation: Near-optimal matrix completion," *IEEE Transactions on Information Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.
- [14] J. F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [15] M. Ashraphijuo, R. Madani, and J. Lavaei, "Characterization of rank-constrained feasibility problems via a finite number of convex programs," in *IEEE 55th Conference on Decision and Control (CDC)*, 2016, pp. 6544–6550.
- [16] M. Ashraphijuo and X. Wang, "Scaled nuclear norm minimization for low-rank tensor completion," *arXiv preprint:1707.07976*, 2017.
- [17] M. Ashraphijuo, R. Madani, and J. Lavaei, "Inverse function theorem for polynomial equations using semidefinite programming," in *IEEE 54th Annual Conference on Decision and Control (CDC)*, 2015, pp. 6589–6596.
- [18] B. Recht and C. Ré, "Parallel stochastic gradient algorithms for large-scale matrix completion," *Mathematical Programming Computation*, vol. 5, no. 2, pp. 201–226, 2013.
- [19] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion using alternating minimization," in *Annual Symposium on the Theory of Computing*, 2013, pp. 665–674.
- [20] R. Ge, J. D. Lee, and T. Ma, "Matrix completion has no spurious local minimum," *arXiv preprint:1605.07272*, 2016.
- [21] D. Pimentel-Alarcón, L. Balzano, and R. Nowak, "Necessary and sufficient conditions for sketched subspace clustering," in *Allerton Conference on Communication, Control, and Computing*, 2016.
- [22] D. Pimentel-Alarcón, N. Boston, and R. D. Nowak, "Deterministic conditions for subspace identifiability from incomplete sampling," in *IEEE International Symposium on Information Theory (ISIT)*, 2015, pp. 2191–2195.
- [23] D. Pimentel-Alarcón, L. Balzano, R. Marcia, R. Nowak, and R. Willett, "Group-sparse subspace clustering with missing data," in *IEEE Workshop on Statistical Signal Processing (SSP)*, 2016, pp. 1–5.
- [24] D. Pimentel-Alarcón, E. R. D. Nowak, and W. EDU, "The information-theoretic requirements of subspace clustering with missing data," in *International Conference on Machine Learning*, 2016.
- [25] D. Pimentel-Alarcón, N. Boston, and R. Nowak, "A characterization of deterministic sampling patterns for low-rank matrix completion," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 623–636, 2016.
- [26] M. Ashraphijuo, X. Wang, and V. Aggarwal, "Deterministic and probabilistic conditions for finite completability of low-rank multi-view data," *arXiv preprint:1701.00737*, 2017.
- [27] M. Ashraphijuo and X. Wang, "Fundamental conditions for low-cp-rank tensor completion," *Journal of Machine Learning Research*, vol. 18 (63), pp. 1–29, 2017.
- [28] M. Ashraphijuo, V. Aggarwal, and X. Wang, "Deterministic and probabilistic conditions for finite completability of low rank tensor," *arXiv preprint:1612.01597*, 2016.
- [29] M. Ashraphijuo and X. Wang, "Characterization of deterministic and probabilistic sampling patterns for finite completability of low tensor-train rank tensor," *arXiv preprint:1703.07698*, 2017.
- [30] M. Ashraphijuo, V. Aggarwal, and X. Wang, "A characterization of sampling patterns for low-Tucker-rank tensor completion problem," in *IEEE International Symposium on Information Theory (ISIT)*, 2017, pp. 531–535.
- [31] M. Ashraphijuo, X. Wang, and V. Aggarwal, "A characterization of sampling patterns for low-rank multi-view data completion problem," in *IEEE International Symposium on Information Theory (ISIT)*, 2017, pp. 1147–1151.
- [32] D. Pimentel-Alarcón and R. Nowak, "A converse to low-rank matrix completion," in *IEEE International Symposium on Information Theory (ISIT)*, 2016.
- [33] M. Ashraphijuo, X. Wang, and V. Aggarwal, "Rank determination for low-rank data completion," *arXiv preprint:1707.00622*, 2017.
- [34] F. Király, L. Theran, and R. Tomioka, "The algebraic combinatorial approach for low-rank matrix completion," *Journal of Machine Learning Research*, vol. 16, pp. 1391–1436, 2015.
- [35] B. Sturmfels, *Solving Systems of Polynomial Equations*. American Mathematical Society, 2002, no. 97.
- [36] E. Candès and B. Recht, "Exact matrix completion via convex optimization," *Communications of the ACM*, vol. 55, no. 6, pp. 111–119, 2012.
- [37] G. Shabat, "Matrix completion using nuclear norm, spectral norm or weighted nuclear norm minimization," 2015.