One-Sided Matrix Completion from Ultra-Sparse Samples

Hongyang R. Zhang Northeastern University, Boston

iter incuster in entiter stug, Beste

Zhenshuo Zhang Northeastern University, Boston

Huy L. Nguyen Northeastern University, Boston

Guanghui Lan

Georgia Institute of Technology, Atlanta

ho.zhang@northeastern.edu

zhang.zhens@northeastern.edu

hu.nguyen@northeastern.edu

george.lan@isye.gatech.edu

Abstract

Matrix completion is a classical problem that has received recurring interest from a wide range of fields. In this paper, we revisit matrix completion in an ultra-sparse sampling setting, where each entry of an unknown, n by d matrix M, is observed with probability $p = \frac{C}{d}$, for any constant $C \ge 2$ (assuming $n \ge d$). This setting is motivated by large-scale panel datasets with high sparsity in practice. While the total number of observed samples, or roughly Cn, is insufficient to recover M, we show that it is possible to recover one side of M, i.e., the second-moment of the row vectors, given by $T = \frac{1}{n}M^{\top}M$. The empirical second moment computed from observational data involves non-random missingness and high sparsity. We design an algorithm that estimates T by normalizing every nonzero entry of the empirical second moment with its observed frequency, followed by gradient descent to impute the missing entries. The normalized entry divides a weighted sum of n binomial random variables by the total number of ones, which is challenging to analyze due to nonlinearity and sparsity. We provide estimation and recovery guarantees for this estimator in the ultra-sparse regime, showing that it is unbiased for any p, and incurs low variance. Assuming the row vectors of M are sampled from a rank-r factor model, we prove that when $n \ge O(\frac{dr^5 \log d}{C^2 \epsilon^2})$, our algorithm can recover T with Frobenius norm error less than ϵ^2 , assuming the rank-r factor model satisfies a standard incoherence condition. We also extend the use of one-sided matrix completion as a sub-procedure towards imputing the missing entries of M.

Experiments on both synthetic and real-world data are provided to evaluate this approach. When tested on three MovieLens datasets, our approach reduces bias by 88% relative to its alternatives. We also validate the linear sampling complexity of n relative to d on synthetic data. On an Amazon reviews dataset with sparsity 10^{-7} , our approach reduces the recovery error of T by 59% and M by 38% compared to existing matrix completion methods.

1 Introduction

Matrix completion (MC) is a classical problem that has received recurring interest from many areas, including compressive sensing (Candes & Plan, 2010), machine learning (Srebro et al., 2004), and signal processing (Chi et al., 2019). A seminal result from the MC literature is that, assuming the unknown matrix M is of low rank and satisfies a certain incoherence assumption — meaning that the row norms of the low-rank factor subspaces are balanced — then exact recovery guarantees from a uniformly random subset of entries are achievable (Candes & Recht, 2012). In this paper, we revisit matrix completion in an ultra-sparse sampling regime, when the sampling probability p is only $\frac{C}{d}$ (for some constant $C \geq 2$), leading to less than $nr \log d$ observed entries — the number of bits needed to represent the underlying low-rank factors of M (assuming

that $n \ge d$). Despite the fact that recovering M is information-theoretically not possible (Chen, 2015), we explore the recovery of one side of M, represented by the second-moment matrix $T = \frac{1}{n}M^{\top}M$.

There are several considerations for revisiting matrix completion in an ultra-sparse sampling regime. First, many large-scale panel datasets are sparse in practice. Consider a recent Amazon reviews dataset, for example (Hou et al., 2024). There are roughly 54 million users times 48 million items, whereas the number of observed entries is about 571 million. Hence, the sparsity ratio is at the order of 10^{-7} . Therefore, we need matrix completion methods that can still perform robustly at this range of sparsity. Second, many recommendation systems require protecting data privacy (Hardt & Roth, 2012). One protocol in this area is that each end user contributes a small amount of data to a central server, which is perturbed with white noise before being collected by the central server (Liu et al., 2015). After receiving data from all users, the central server computes a summary statistic, such as a low-rank subspace spanned by row vectors, which can be publicly released. Each user then uses this summary statistic locally by projecting their data onto a low-rank subspace (Wang et al., 2023).

Recent work (Cao et al., 2023) initiates the study of one-sided matrix completion, whose objective is to estimate the averaged second moment matrix of the row vectors of M (i.e., T), in the setting where two entries from each row are observed. Their work leaves open the general case when more than two entries are observed in each row. To illustrate, let \hat{M} denote the observed data matrix and I denote the zero-one mask of \hat{M} . As mentioned earlier, the observed entries from \hat{M} are insufficient to recover M, rendering existing MC results not applicable to one-sided completion. Let Ω denote the nonzero index set of the empirical second moment $\hat{M}^{\top}\hat{M}$. First, the missingness in Ω is non-random because the diagonal entries are fully observed (with high probability), while the off-diagonal entries are only sparsely observed (See Claim 2.1). Further, the observation frequencies in the off-diagonal entries, as measured by $I^{\top}I$, are small in the ultra-sparse regime. Second, while stochastic gradient algorithms are widely used in practice (Lan, 2020), theoretically analyzing their performance for MC is challenging (Sun & Luo, 2016) and has not been done when Ω involves non-random missingness.

To address the first issue above, we consider a weighted estimator, which normalizes every nonzero entry of $\hat{M}^{\top}\hat{M}$ by the corresponding entry at $I^{\top}I$. This method of normalizing a sum of weighted binomial random variables has been studied before in survey sampling (Särndal et al., 2003) and is also known as Hájek's estimator (Hájek, 1971) (See, e.g., Hirano et al. (2003) from the causal literature). However, theoretically analyzing Hájek's estimator is challenging because it involves dividing two random variables, which is a nonlinear operation. Our first key observation is that Hájek's estimator is unbiased on Ω , the observed subset of T. Moreover, by applying a first-order approximation on the variance of the normalized entry of $\hat{M}^{\top}\hat{M}$, we find that Hájek's estimator achieves low variance (relative to the Horvitz-Thompson estimator (Horvitz & Thompson, 1952), see Section 2). The algorithmic implication is that Hájek's estimator provides more stable performance, which leads to low bias on Ω (as we will verify in experiments).

As for the second issue, we provide recovery guarantees for any local minimizer of a loss objective between the normalized second-moment matrix and a low-rank factorization, for imputing missing entries of T that are outside Ω . Under a rank-r factor model on the row vectors of M, we show that when $n \geq O\left(\frac{dr^5 \log d}{C^2 \epsilon^2}\right)$, then any local minimizer of the loss objective is within ϵ^2 Frobenius norm distance to T. This result is established assuming the rank-r factor model satisfies an incoherence assumption — Otherwise, one can find vacuous examples where recovering T is not possible (Similar to the ones in Candes & Recht (2012)). As a remark, the bound on n relative d is optimal. Further, our result applies to any constant C, showing that one-sided MC is possible in the general case, thus resolving an open question from Cao et al. (2023). A key technical result is establishing a concentration inequality to tackle the non-random missingness of Ω , by leveraging a conditional independence property at each row of Ω . This conditional independence property allows us to bound the bias of Hajek's estimator row-by-row.

We extensively validate our approach on both synthetic data and real-world datasets. We find that Hájek's estimator is particularly effective in sparse sampling regimes. For instance, when each row has two observed entries, our approach can outperform a prior nuclear norm regularization method by 70% in terms of recovering T. Additionally, we explore the use of one-sided matrix completion as a sub-procedure for imputing missing entries of M. On the Amazon reviews dataset, we extend our algorithm to impute M via a least squares

Table 1: Comparison between this work and several closely-related prior results on matrix completion. Here M denotes an n by d matrix with $n \ge d$, and $T = \frac{1}{n}M^{\top}M$ is the averaged second-moment of the row vectors. In particular, we focus the comparison on the one-sided matrix completion or ultra-sparse sampling settings.

Estimand	Sampling Regime	Main Approach and Reference
M	$O(\frac{1}{n})$	Thresholded Alternating Minimization (Gamarnik et al., 2017)
T	Two entries per row	Nuclear Norm Regularization (Cao et al., 2023)
T	$p = \frac{C}{n} \ (\forall C \ge 2)$	HÁJEK-GD with Incoherence Regularization (<i>This paper</i>)

regression after estimating T. We find that this approach can outperform alternating gradient descent and soft impute with alternating least squares (Hastie et al., 2015) — both widely used methods for matrix completion (Chi et al., 2019; Li et al., 2020) — by 21% and 38% in terms of the root mean squared error, for recovering the missing entries of M. Ablation studies show that Hájek's estimator reduces bias on Ω by 99% relative to Horvitz-Thompson's estimator on synthetic data and by 88% averaged over three MovieLens datasets. Incorporating an incoherence regularization penalty into the loss (Ge et al., 2016) reduces the recovery error of T by 22% on synthetic data and by 52% on the Amazon reviews dataset.

Summary of contributions. This paper revisits the matrix completion (MC) problem in an ultra-sparse sampling regime when there are only Cn samples from an unknown n by d matrix, for some constant $C \geq 2$. Our contributions to the MC literature are three-fold: (1) Designing a new algorithm for one-sided matrix completion via Hájek's estimator, along with a thorough analysis of Hájek's estimator in MC, showing that it incurs lower variance than Horvitz-Thompson's estimator. (2) Providing near-optimal sample complexity bounds on n relative to d for one-sided matrix completion in a low-rank model, which now applies to gradient-based optimization and more general sampling regimes than existing results. See Table 1 for a summary of comparisons between our results and several related works. (3) Experiments validating the effectiveness of our approach on sparse matrix datasets, including new results of using one-sided matrix completion for full matrix completion. The code for replicating these experiments can be accessed at: https://anonymous.4open.science/r/One-sided-matrix-completion-ultra-sparse-samples.

Organizations. The rest of this paper is organized as follows. In Section 2, we describe the problem setup and several standard notations. In Section 3, we present our approach and its theoretical analysis. We present a proof sketch in Section 4. Then, we present experiments in Section 5. We discuss related works in Section 6. Lastly, we summarize the paper in Section 7. We present complete proofs in Appendix A and B, and omitted experiments in Appendix C.

2 Preliminaries and Notations

Let M denote an n by d, unknown matrix. Assume the rank of M is equal to r. With loss of generality, suppose that $n \ge d$. For example, each row of M may include the ratings of a user, while each column may correspond to the ratings of an item. We observe a partial matrix, denoted by \hat{M} . Let $I \in \{0,1\}^{n \times d}$ denote the indicator matrix on the observed entries of \hat{M} . Our goal is to recover the averaged second-moment matrix $T := n^{-1}M^{\top}M$, given \hat{M} and I.

We assume that each entry of \hat{M} is observed independently with probability $p \in (0, 1)$, following prior literature (e.g., Candes & Recht (2012); Sun & Luo (2016); Ge et al. (2016)). We focus on a constant sampling regime, where $p = \frac{C}{d}$, for some fixed integer $C \ge 2$. In such a regime, m is roughly Cn, less than the number of bits to represent an n by r matrix.

An important quantity for working with one-sided estimation is $I^{\top}I$, which is a symmetric matrix in dimension d. This matrix provides the empirical frequency counts for the empirical second moment matrix $\hat{M}^{\top}\hat{M}$. For example, consider an off-diagonal entry (i, j), where $i \neq j$ and both i, j are in [d]. The (i, j)-th entry of $I^{\top}I$, denoted as $(I^{\top}I)_{i,j}$, is equal to the number of overlapping, nonzero entries between the *i*-th and *j*-th column

of \hat{M} :

$$(I^{\top}I)_{i,j} = \sum_{k=1}^{n} I_{k,i} I_{k,j}.$$
(1)

Let Ω denote the indices of the nonzero entries of this matrix and I denote the zero-one indicator mask of Ω . We now make the following observation regarding the observed entries in Ω . In particular, we find that the sampling patterns of Ω are non-random in the following sense.

Claim 2.1. For every i = 1, 2, ..., d, the diagonal entries satisfy that

$$\Pr[(i, i) \in \Omega] = 1 - (1 - p)^n.$$

The off-diagonal entries satisfy that for every i = 1, ..., d, and j = 1, ..., d such that $j \neq i$,

$$\Pr[(i, j) \in \Omega] = 1 - (1 - p^2)^n.$$

To illustrate, recall that p = C/d. Suppose $n = O(d \log d)$. The probability that $(i, i) \in \Omega$ is roughly $1 - e^{-pn} \approx 1 - d^{-O(C)}$. Then by a union bound, with probability at least $1 - O(d^{-1})$, it must be the case that $(i, i) \in \Omega$, for every i = 1, 2, ..., d.

For each off-diagonal entry, it is observed in Ω with probability equal to $1 - (1 - p^2)^n$. For simplicity of notation, let us denote this as q in the rest of the paper. Note that q roughly equal to $np^2 = O(\frac{C^2 \log d}{d})$.

To account for the difference in sampling probabilities between diagonal and off-diagonal entries of Ω , let $P_{\Omega} : \mathbb{R}^{d \times d} \to \mathbb{R}^{d \times d}$ be a weighted projection operator defined as follows:

$$(P_{\Omega}(Z))_{i,j} = \begin{cases} qZ_{i,i} & \text{if } (i,i) \in \Omega, \\ Z_{i,j} & \text{if } (i,j) \in \Omega \text{ and } j \neq i, \\ 0 & \text{if } (i,j) \notin \Omega. \end{cases}$$

Having introduced the mask matrix $I^{\top}I$ and Ω , next, we discuss the normalization of the empirical second moment of the observed matrix, given by $\hat{M}^{\top}\hat{M}$. One way to normalize $\hat{M}^{\top}\hat{M}$ for one-sided matrix completion is via the Horvitz-Thompson (HT) estimator (Horvitz & Thompson, 1952), which inversely reweights each entry with the true sampling probability (i.e., p on the diagonal entries and p^2 on the off-diagonal entries):

$$\overline{T}_{i,j} = \begin{cases} \frac{\sum_{k=1}^{n} M_{k,i}^{2} I_{k,i}}{np}, & \text{if } i = j, \\ \frac{\sum_{k=1}^{n} M_{k,i} M_{k,j} I_{k,i} I_{k,j}}{np^{2}}, & \text{if } i \neq j. \end{cases}$$

One can verify that $\mathbb{E}[\overline{T}_{i,j}] = T_{i,j}$. Notice that np^2 is the expectation of $(I^{\top}I)_{i,j}$ (cf. equation (1)), while np is the expectation when i = j. When p is too small, np^2 becomes small, resulting in high fluctuation in the off-diagonal entries. Another approach, known as Hajek's estimator (Hájek, 1971), is to normalize each entry in Ω by the estimated sampling probability, which is canceled out since the probability value is the same across different k, leading to the following expression, for any $(i, j) \in \Omega$:

$$\widehat{T}_{i,j} = \begin{cases} \frac{\sum_{k=1}^{n} M_{k,i}^{2} I_{k,i}}{\sum_{k=1}^{n} I_{k,i}}, & \text{if } i = j, \\ \frac{\sum_{k=1}^{n} M_{k,i} M_{k,j} I_{k,i} I_{k,j}}{\sum_{k=1}^{n} I_{k,i} I_{k,j}}, & \text{if } i \neq j. \end{cases}$$

The theoretical analysis of Hájek's estimator is challenging due to its nonlinear form, which involves dividing one random variable by another. It is known that Hájek's estimator is asymptotically unbiased (Hájek, 1971) in the large *n* limit, incurring an error of $O(n^{-1/2})$. In addition, Hájek's estimator provides a variance reduction effect compared to Horvitz-Thompson's estimator in causal inference (Hirano et al., 2003). Not much is known in the MC context, and several natural questions arise. First, how does this estimator work for one-sided matrix completion? What is the bias of \hat{T} , and how does Hájek's estimator compare to Horvitz-Thompson's estimator? Second, how can we use \hat{T} to impute the missing entries outside Ω ? The rest of this paper is dedicated to answering these questions. **Notations.** Before continuing, we provide a list of notations for describing the results. Following the convention of big-O notations, given two functions f(n) and g(n), let $f(n) \leq O(g(n))$ indicate that there exists a constant c independent of n such that when $n \geq n_0$ for some large enough n_0 , then $f(n) \leq c \cdot g(n)$. We use the notation $f(n) \leq g(n)$ as a shorthand for indicating that f(n) = O(g(n)). Let $\tilde{O}(g(n))$ denote that $O(\log^c(n)g(n))$ for some constant c independent of n. We also use the little-o notation $f(n) \leq o(g(n))$, which means that f(n)/g(n) goes to zero as n goes to infinity.

Let $[d] = \{1, 2, ..., d\}$ denote a shorthand notation for the set from 1 up to d. We use $\mathcal{N}(a, b)$ to denote a Gaussian distribution with mean set at a and variance equal to b. We use the notation $(x)_{+} = \max(x, 0)$ to denote a truncation operation set above zero.

Let $\|\cdot\|_F$ denote the Frobenius norm of an input matrix. Let $\|\cdot\|_2$ denote the spectral norm of the input matrix and let $\|\cdot\|$ denote the Euclidean norm of a vector. Let $\langle X, Y \rangle = \text{Tr}[X^\top Y]$ denote the inner product between two matrices that have the same dimensions. Let $\|\cdot\|_{\infty}$ denote the infinity norm of a matrix, which corresponds to its largest entry in absolute value.

3 Estimation and Recovery Guarantees

In this section, we present a thorough analysis of Hájek's estimator in the one-sided matrix completion problem. We begin by analyzing the bias of \hat{T} compared to T on the observed entries Ω . A key observation is that \hat{T} provides an unbiased estimate of T, as stated formally below.

Lemma 3.1. Suppose the entries of \hat{M} are sampled from M independently with a fixed probability $p \in (0,1)$.¹ Let Ω denote the index set corresponding to the non-zero entries of $\hat{M}^{\top}\hat{M}$. Then, the following must be true:

$$\mathbb{E}\left[\widehat{T}_{i,j}\middle|(i,j)\in\Omega\right] = T_{i,j},\tag{2}$$

for any $1 \leq i, j \leq d$. In other words,

$$\mathbb{E}\left[\widehat{T}_{i,j}\right] = \Pr\left[(i,j) \in \Omega\right] \cdot T_{i,j}.$$

Unlike conventional results that show Hajek's estimator is asymptotically consistent, in our setting, we find that it is unbiased in the finite sample regime. We describe the main ideas for showing Lemma 3.1. We note that for any k > 1, conditioned on $(I^{\top}I)_{i,j} = k$, $\hat{M}^{\top}\hat{M}$ must be equal to the sum of k pairs chosen from $M_{1,i}M_{1,j}, M_{2,i}M_{2,j}, \ldots, M_{n,i}M_{n,j}$. Furthermore, the choice of k such pairs out of n is uniformly at random among all possible $\binom{n}{k}$ combinations, because the choices are symmetric. Thus, the average of k randomly chosen pairs must be equal to the average of all n pairs in expectation. This is stated precisely as follows:

$$\mathbb{E}\left[\frac{1}{k}\left(\hat{M}^{\top}\hat{M}\right)_{i,j}\middle|\left(I^{\top}I\right)_{i,j}=k\right]=\frac{1}{n}\sum_{a=1}^{n}M_{a,i}M_{a,j}=T_{i,j}.$$
(3)

3.1 Variance Reduction

Having established that \widehat{T} provides an unbiased estimate of T, we study the variance of \widehat{T} . We show that Hj́ake's estimator incurs lower variance relative to the Horvitz-Thompson estimator. We first derive an approximation of $Var(\widehat{T})$ to tackle the nonlinearity of Hájek's estimator.

Theorem 3.2. Suppose $p = \frac{C}{d}$ for some fixed integer $C \ge 2$. Suppose $n \ge O(d \log d)$. Then, with probability at least $1 - O(d^{-1})$, all of the diagonal entries of \widehat{T} are in Ω , and for any diagonal entries of \widehat{T} , the variance of $\widehat{T}_{i,i}$, for any i = 1, 2, ..., d, is approximated by

$$Var\left(\widehat{T}_{i,i}\right) = \frac{1-p}{np} \left(\frac{1}{n} \sum_{k=1}^{n} M_{k,i}^{4}\right) - \frac{1-p}{np} \cdot T_{i,i}^{2} + O\left(\frac{1}{np}\right).$$

$$\tag{4}$$

¹We remark that the same unbiased estimation result can also be stated in the setting where every row has k uniformly random entries from M, for arbitrary integer $k \ge 2$. The proof is similar to that of Lemma 3.1.



Figure 1: An illustration of the variance reduction using Hájek's estimator vs. Horvitz-Thompson's estimator, measured on synthetic data with $n = 10^4$ and $d = 10^3$, by repeating the data sampling procedure 100 times and calculating the variance across either uniform sampling with probability p, or sampling C entries per row. In particular, $Var(\hat{T})$ is generally 10^{-3} lower than $Var(\bar{T})$. For details on the simulation setup, see Section 5.

For any nonzero, off-diagonal entries in Ω , the variance of $\widehat{T}_{i,j}$, for any $1 \leq i \neq j \leq d$, is approximated by

$$Var\left(\widehat{T}_{i,j}\Big|(i,j)\in\Omega\right) = \frac{1}{n}\left(\sum_{k=1}^{n} M_{k,i}^2 M_{k,j}^2\right) - T_{i,j}^2 + O\left(\frac{(\log d)^2}{d}\right).$$
(5)

We now compare $Var(\widehat{T})$ with $Var(\overline{T})$, since both estimators are unbiased. We first derive the variance of \overline{T} to make the comparison. For diagonal entries, we have that $\overline{T}_{i,i} = (np)^{-1}(\hat{M}^{\top}\hat{M})_{i,i}$, for $i = 1, \ldots, d$. For off-diagonal entries, we have $\overline{T}_{i,j} = (np^2)^{-1}(\hat{M}^{\top}\hat{M})_{i,j}$, for $i \neq j$. Then, we calculate variance as:

$$Var(\overline{T}_{i,i}) = \frac{1-p}{n^2 p} \sum_{k=1}^{n} M_{k,i}^2, \ \forall i = 1, \dots, d,$$
(6)

$$Var\left(\overline{T}_{i,j}\right) = \frac{1-p^2}{n^2 p^2} \sum_{k=1}^n M_{k,i}^2 M_{k,j}^2, \ \forall j = 1, \dots, d, \text{ and } j \neq i.$$
(7)

Now we can compare the above with equation (4) and (5), respectively. Notice that the diagonal entries are nonzero with high probability over all the diagonal entries. By comparing equation (4) to equation (6), we see a reduction in the second term from equation (4), which is at the same order as the leading term.

The off-diagonal entries are nonzero with probability equal to $q \approx np^2 = \tilde{O}(d^{-1})$. By comparing equation (5) to (7), we see that equation (7) is larger than the first term of equation (5) by a factor of np^2 . In addition, the second term of equation (5), which is always negative, further contributes to the reduced variance of $\hat{T}_{i,j}$ compared with $\overline{T}_{i,j}$.

Algorithmic implications. One consequence of variance reduction is that the empirical performance of Hájek's estimator tends to be much more stable than Horvitz-Thompson in the ultra-sparse setting. This is illustrated in Figure 1, measured on synthetic datasets with $n = 10^4$ and $d = 10^d$, for various values of p (uniform sampling) and C (i.e., sampling C entries per row without repetition).

3.2 Sample Complexity

Next, we consider the entries outside Ω . We minimize a reconstruction error objective between \widehat{T} and XX^{\top} , for some $X \in \mathbb{R}^{d \times r}$, plus a regularization penalty of R(X), as follows:

$$\ell(X) := \frac{1}{2} \left\| P_{\Omega}(XX^{\top} - \widehat{T}) \right\|_{F}^{2} + \lambda R(X),$$
(8)

where $\lambda > 0$ is a regularization parameter. R(X) is defined as $\sum_{i=1}^{d} (||X_i|| - \alpha)_+^4$, where $X_i \in \mathbb{R}^r$ is the *i*-th row vector of X, for $i = 1, \ldots, d$, and α is a scalar that roughly corresponds to the row vector norms of X. This penalty regularizes the maximum ℓ_2 -row-norm of X, and is shown to provide theoretical properties in

Algorithm 1 Hájek estimation with gradient descent (HÁJEK-GD) for one-sided matrix completion

Input: A partially-observed data matrix $\hat{M} \in \mathbb{R}^{n \times d}$

Require: Rank of the second-moment matrix r, number of iterations t, and learning rate η **Output:** A d by d matrix

1: $I \in \mathbb{R}^{n \times d} \leftarrow$ The 0-1 indicator mask corresponding to the nonzero entries of \hat{M}

2: $\Omega \subseteq [d] \times [d] \leftarrow$ The set of indices corresponding to the nonzero entries of $I^{\top}I$

- 3: $\hat{T} \in \mathbb{R}^{d \times d} \leftarrow$ The element-wise division between $\hat{M}^{\top} \hat{M}$ and $I^{\top} I$ on Ω
- 4: $X_0 \in \mathbb{R}^{d \times r} \leftarrow A$ random Gaussian matrix whose entries are sampled independently from $\mathcal{N}(0, d^{-1})$
- 5: **for** i = 1, ..., t **do**
- 6: $X_i \leftarrow X_{i-1} \eta \nabla \ell(X_{i-1})$, where $\ell(\cdot)$ is defined in equation (8)
- 7: end for
- 8: **return** \widehat{T} , plus the entries of $X_t X_t^{\top}$ outside Ω

the optimization landscape of matrix completion (Ge et al., 2016; 2017). We summarize our procedure in Algorithm 1, which uses gradient descent to minimize $\ell(X)$.

We analyze this algorithm in a common means model, where each row of M follows a mixture of r vectors $u_1, u_2, \ldots, u_r \in \mathbb{R}^d$. Let M_i denote the *i*-th row vector of M. For every $i = 1, 2, \ldots, n$, assume that M_i is drawn uniformly at random from u_1, u_2, \ldots, u_r .² Let $U = [u_1, u_2, \ldots, u_r]/\sqrt{r}$ denote a d by r matrix corresponding to the combined rank-r factors. Let $\kappa = \sigma_{\max}(U)/\sigma_{\min}(U)$ denote the condition number of U. A critical condition to ensure guaranteed recovery in matrix completion is the following assumption regarding the row norms of U.

Assumption 3.3 (See also Definition 1.1, Recht (2011) and Assumption 1, Ge et al. (2016)). Let Ue_i denote the *i*-th row vector of *a d* by *r* matrix *U*, where e_i is the *i*-th basis vector, for i = 1, 2, ..., d. The coherence of *U* is given by

$$\mu := \frac{d}{r} \max_{1 \le i \le d} \frac{\|Ue_i\|^2}{\|U\|_F^2},\tag{9}$$

Assuming that $\mu(U)$ is a fixed value that does not row with d, we show the following performance guarantee of gradient descent for recovering T, measured in terms of the Frobenius norm distance between a local minimizer and T.

Theorem 3.4. Suppose the rows of M follow a mixture of r common factors given by $U \in \mathbb{R}^{d \times r}$. Additionally, the coherence of U is at most μ for some fixed $\mu \geq 1$ that does not grow with d. Suppose each entry of M is observed with probability $p = \frac{C}{d}$ for some fixed integer $C \geq 2$, and let $q = 1 - (1 - p^2)^n$. Let $\alpha = 4\kappa^2 r \sqrt{\frac{\mu}{d}}$ and $\lambda = \frac{(r+1)dq}{16r^2\mu^3}$. When $n \geq \frac{cdr^5\kappa^6\mu^2\log(d)}{C^2\epsilon^2}$ for some fixed constant c and $\epsilon \in (0,1)$, with probability at least $1 - O(d^{-1})$ over the randomness of Ω , when d is large enough, any local minimizer X of $\ell(\cdot)$ satisfies

$$\left\|XX^{\top} - T\right\|_{F}^{2} \le \epsilon^{2}.$$
(10)

This result implies that any local minimum solution of the loss objective is also approximately a global minimum solution. We defer a proof sketch of this result to Section 4. We remark that prior results have established one-sided matrix completion guarantees when each row has two observed entries (Cao et al., 2023). By contrast, Theorem 3.4 applies to arbitrary values of $C \ge 2$ in each row.

Algorithmic implications. A crucial step in the proof is to leverage the incoherence assumption, which reduces local optimality conditions from finite samples to the entire population. Later in Section 5.3, we also empirically demonstrate that the use of the incoherence regularization penalty helps improve performance on real-world datasets.

²As a remark, this sort of "common means" model is commonly used to study heterogeneous data. See, e.g., Kolar et al. (2011); Dobriban & Sheng (2020). It is worth mentioning that the proof can be easily extended to handle noise addition, such as a noise vector ϵ_i added to M_i , whose entries are drawn independently from a Gaussian distribution $\mathcal{N}(0, \sigma^2/d)$.

4 Proof Techniques and Extensions

Next, we present a sketch of our proof. There are two major challenges in analyzing Hájek's estimator in the ultra-sparse sampling regime: First, the estimator involves the division of two random variables, resulting in a nonlinear estimator. Further, the missing patterns in Ω are non-random, depending on the diagonal and off-diagonal entries, respectively. Our approach to tackle this non-linearity is via a first-order approximation technique for analyzing Hájek's estimator on the diagonal entries. Further, we carefully analyze the bias in the off-diagonal entries.

Second, the sampling patterns of Ω are not fully independent. To this end, we establish a concentration inequality that handles the concentration error row by row in Ω . This is based on the observation that, conditioned on one row, expect the diagonal entries, the randomness of each entry in that row would be independent across different entries. This row-by-row concentration argument also allows us to analyze the spectral norm of the bias matrix. As a remark, this type of argument has been in matrix completion with non-random missing data (Athey et al., 2021). However, prior work focuses on analyzing nuclear norm regularization methods, while our result now applies to gradient-based optimization. Next, we outline the main results in addressing the above two challenges.

Bias of Hájek's estimator. A key result for deriving the variance, or expected squared error of Hájek's estimator, since \hat{T} is unbiased by Lemma 3.1, is as follows.

Lemma 4.1. In the setting of Theorem 3.4, suppose p = C/n for some fixed integer $C \ge 2$ and $n \ge d(\log d)/\epsilon^2$. With probability at least $1 - O(d^{-1})$, the following must be true, for every i = 1, 2, ..., d,

$$Var\left(\widehat{T}_{i,i}\right) \leq \frac{\mu^2 r}{d^2 n p} \left(1 + \epsilon \sqrt{\frac{6C^{-1}}{d}}\right) \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}}\right),\tag{11}$$

$$\sum_{i \in S_i} Var\left(\widehat{T}_{i,j}\right) \le \frac{\mu q}{d} \left(1 + 3r\sqrt{\frac{3\log(d/2)}{2d}}\right) + \widetilde{O}\left(d^{-3}\right).$$

$$\tag{12}$$

Proof sketch. The proof regarding diagonal entries in Eq. (11) is based on a first-order approximation of $Var(\hat{T}_{i,i})$. For every i = 1, 2, ..., d, let

$$A_{i,i} := \sum_{k=1}^{n} M_{k,i}^2 I_{k,i} \text{ and } B_{i,i} := \sum_{k=1}^{n} I_{k,i}.$$
(13)

We perform Taylor's expansion around $(\mathbb{E}[A_{i,i}], \mathbb{E}[B_{i,i}])$ as follows (See also Chapter 5.5, Särndal et al. (2003)):

$$\widehat{T}_{i,i} = \frac{\mathbb{E}\left[A_{i,i}\right]}{\mathbb{E}\left[B_{i,i}\right]} + \frac{1}{\mathbb{E}\left[B_{i,i}\right]} (A_{i,i} - \mathbb{E}\left[A_{i,i}\right]) - \frac{\mathbb{E}\left[A_{i,i}\right]}{(\mathbb{E}\left[B_{i,i}\right])^2} (B_{i,i} - \mathbb{E}\left[B_{i,i}\right]) + \epsilon_{i,i}.$$
(14)

Above, the first term to the right of equation (14) stems from the zeroth-order expansion. The second and third terms arise from taking the partial derivatives over $A_{i,i}, B_{i,i}$, respectively. $\epsilon_{i,i}$ is an error term that is at the order of the variance of $A_{i,i}, B_{i,i}$, which can be calculated in close form based on equation (13).

Given that \hat{T} is unbiased, $Var(\hat{T}_{i,i})$ is equal to the expectation after squaring both sides of Eq. (14). As a result, we derive the variance approximation of $\hat{T}_{i,i}$ by carefully analyzing each term, which will lead to equation (4). Figure 2 provides an illustration of equations (14) and (4) except the $\epsilon_{i,i}$ error term. We find that both approximations hold with negligible errors for various sampling probabilities p (uniform sampling with probability p) and C (fixed number of entries per row). This simulation uses $n = 10^4$ and $d = 10^3$ and follows the setup stated in Section 5.

For off-diagonal entries, note that when $(i, j) \in \Omega$, with high probability the (i, j)-th entry of $I^{\top}I$ must be one. Thus, $Var(\hat{T}_{i,j})$ is equal to the variance of n values $M_{1,i}M_{1,j}, \ldots, M_{n,i}M_{n,j}$ (plus some small errors). This leads to the variance approximation of equation (5).



(a) Taylor's expansion (b) Taylor's expansion (c) Variance approximation (d) Variance approximation

Figure 2: We illustrate the results from applying a first-order approximation to the variance of the diagonal entries of HÁJEK-GD, run on synthetic data with $n = 10^4$ and $d = 10^3$. Figures 2a and 2c: Sample each entry with probability p. Figures 2b and 2d: Sample C entries per row without repetition. In particular, the approximation errors incurred by both first-order Taylor's expansion and the variance approximation are generally less than 10^{-6} .

A concentration inequality for P_{Ω} . In order to analyze the optimization landscape of $\ell(\cdot)$, we will first derive the first- and second-order optimality conditions. However, these are stated on the observed set Ω , and we need to turn them into the full set on T. Towards that end, we develop the following concentration inequality, which helps reduce finite-sample local optimality conditions to the full matrix.

Lemma 4.2. Let $W \in \mathbb{R}^{d \times d}$ be a d by d symmetric matrix such that $||W||_{\infty} \leq \frac{\nu}{d}$, for some fixed positive value of ν . Then, with probability at least $1 - O(d^{-1})$, the following statement must be true:

$$\|P_{\Omega}(W) - qW\|_2 \le q\nu \sqrt{\frac{16\log(2d)}{dq}}.$$
 (15)

The proof of this concentration result can be found in Appendix B.2. Finally, we analyze the population loss of one-sided matrix completion, building on the machinery of Ge et al. (2016). We extend their analysis to account for the bias of \hat{T} and the non-random missingness of Ω (cf. Claim 2.1). In particular, the treatment of the bias of \hat{T} relative to T is especially tedious and requires careful calculation. We build on Lemma 4.1 and 4.2 to give the bias of each diagonal entry, and the bias of each row — See Corollary 4.1. A more detailed comparison is presented after we present the full proof in Remark B.11. The complete proof of Theorem 3.4 can be found in Appendix B.

Extensions. We now describe two extensions of our approach. First, to recover individual entries of M using HÁJEK-GD, we solve a least squares problem by projecting the observed entries onto the span of the recovered second moment. See Figure 3 for an illustration of the overall procedure. The full procedure is described in Algorithm 2, Appendix C.³

Second, we can extend the recovery guarantee to account for differential privacy using the Gaussian mechanism. An algorithm $\mathcal{A}(\hat{M})$ is said to satisfy (ε, δ) -joint differential privacy, if, for an arbitrary i = 1, 2, ..., n, after replacing the *i*-th row of M (and the *i*-th row of \hat{M} , correspondingly) by another vector x from domain $\mathcal{X} \subseteq \mathbb{R}^d$, for any set of events S, the output of \mathcal{A} satisfies the following:

$$\Pr\left[\mathcal{A}(\hat{M}) \in S\right] \le \exp(\varepsilon) \Pr\left[\mathcal{A}(\hat{M}') \in S\right] + \delta, \tag{16}$$

where we use \hat{M}' to denote the matrix whose *i*-th row of \hat{M} has been replaced by some $x \in \mathcal{X}$ and the rest of the matrix remains the same as \hat{M} . This notion of joint-differential privacy is based on a line of related work from the differentially-private matrix completion problem (Liu et al., 2015; Wang et al., 2023). By using the

³As mentioned in the introduction, in general, it is not possible to provide recovery guarantees since the number of samples is too few. One way to view this algorithm is that when a user runs the least squares regression locally, it is likely that they have more data, e.g., with O(r) entries, the least squares procedure can accurately recover the entire row, based on an accurately-recovered T.



Figure 3: Let \hat{M} denote an n by d partially observed matrix. For instance, each row of M corresponds to the review ratings of a user, and each column represents information about an item. Thus, M is a tall-and-skinny matrix, with n being larger than d. Let Ω denote the non-zero entries of $\hat{M}^{\top}\hat{M}$, the set of observed entries in $\hat{M}^{\top}\hat{M}$. **Step 1:** The missing patterns are non-random. We tackle this non-random missingness using Hájek's estimator, applied to every entry of $\hat{M}^{\top}\hat{M}$ on Ω . **Step 2:** Impute the remaining entries outside Ω by running gradient descent on a reconstruction error objective between a low-rank factorization and \hat{T} and finally fill in the missing entries outside Ω using $X_t X_t^{\top}$. The estimated T is the sum of \hat{T} plus the imputed entries from $X_t X_t^{\top}$ outside Ω . Let U_r denote a rank-r subspace computed via SVD of the estimated T. One can recover an approximation of M by solving a least squares regression with variables $Q \in \mathbb{R}^{n \times r}$. For details, see Algorithm 2.

Gaussian mechanism to perturb the nonzero entries of \hat{M} with a random sample from $\mathcal{N}(0, \sigma^2)$, we measure the ℓ_2 -sensitivity of an algorithm $\mathcal{A} : \mathbb{R}^{n \times d} \to \mathbb{R}^{d \times d}$ as:

$$\Delta_2(\mathcal{A}) = \max_{M \sim M'} \left\| \mathcal{A}(M) - \mathcal{A}(M') \right\|_F.$$
(17)

For $\sigma = 2\sqrt{\ln(1.25\delta^{-1})}\Delta_2(\mathcal{A})/\varepsilon$, with high probability, HÁJEK-GD satisfies the (ε, δ) -joint differential privacy. See Theorem A.1, Dwork & Roth (2014). In Appendix C.2, we examine the sensitivity of HÁJEK-GD and find that it remains low in practice.

5 Experiments

Lastly, we evaluate our approach through extensive experiments on both synthetic data and real-world datasets. We generate synthetic data by sampling each entry of M from a Gaussian distribution as $\mathcal{N}(1/\sqrt{d}, 1/d)$.⁴ Then, we observe \hat{M} by independently sampling each entry with probability p = C/d.⁵ As for the experiments on real-world data, we consider three large, sparse panel datasets, including MovieLens datasets with 20 million to 32 million nonzero entries, a recent Amazon Reviews dataset (Hou et al., 2024), and another dataset based on the 1,000 Genomes Project (Cao et al., 2023). These matrix datasets are relatively standard, but they represent typical sparse datasets that one might encounter in practice. We refer to further details about the datasets, such as their statistics, in Appendix C.⁶

The baseline methods under comparison include i) nuclear norm regularization (Cao et al., 2023; Zhang et al., 2019b),⁷ ii) alternating gradient descent (alternating-GD), and iii) soft impute with alternating least squares (softImpute-ALS) (Hastie et al., 2015). All three methods are widely studied in the matrix completion

⁴Unless otherwise stated, we fix $n = 10^4$ and $d = 10^3$ for all the synthetic data. We apply SVD to M and obtain the top-10 singular values and their corresponding singular vectors as M; in other words, the lower tail of M is truncated.

⁵Following prior work (Cao et al., 2023), we also experiment with sampling C entries from every row, though we observe that the results are largely the same.

 $^{^{6}}$ We set up all three datasets as regression tasks. The implementation uses PyTorch and runs on an Ubuntu server with 16 Intel Xeon CPUs and an Nvidia Quadro 6,000 GPU card. Further details can be found in Appendix C.

 $^{^{7}}$ We note that their approach weighted estimator as HÁJEK-GD when each row contains two entries. One difference is that we add the incoherence regularization to the loss objective.



(a) Comparing the bias of \widehat{T} with \overline{T}

(b) Comparing the variance of \widehat{T} with the variance of \overline{T}

 \overline{T} : Uniform

 \overline{T} : Non-uniform

Figure 4: Comparing the bias and variance of Hájek's estimator vs. Horvitz-Thompson's estimator, measured as mean squared error between \hat{T} (or \overline{T}) and T on Ω . Left: Synthetic data with $n = 10^4$ and $d = 10^3$. Right: MovieLens datasets. We consider both uniform sampling with probability p and sampling C entries from each row. In Figure 4b, we further consider a non-uniform sampling setting, where a row with more observed entries also has a higher sampling probability.

literature (See, e.g., a recent work for references (Chi et al., 2019)). To use ii) and iii) in the one-sided matrix completion, we first recover M and then compute T with the recovered M.

We compare the estimated second-matrix matrix to the true T by measuring the Frobenius norm of their difference, and refer to this as the recovery error of an algorithm. For recovering the missing entries of M, we compare Algorithm 2 to alternating-GD and softImpute-ALS, and report the root mean square error (RMSE) between the imputed entries and their true values.

Summary of numerical results. For one-sided matrix completion in the most sparse setting where each row has two entries, we find that HÁJEK-GD reduces recovery error by 70% compared to nuclear norm regularization on synthetic data. For all three real-world datasets (whose sparsity level is less than 1%), HÁJEK-GD reduces recovery error of T by at least 42% compared to alternating-GD, and 59% compared to softImpute-ALS.

By using HÁJEK-GD for imputing the missing entries of M in Algorithm 2, we find that this approach reduces RMSE by 85% compared to alternating-GD and softImpute-ALS, on synthetic data where each row has only two entries. For all three real-world datasets, HÁJEK-GD reduces RMSE by at least 21% compared to alternating-GD, and 38% compared to softImpute-ALS.

Finally, we report ablation studies to validate our algorithmic insights, including the benefits of adding the regularization R(X) to the loss and the low sensitivity of HÁJEK-GD to random noise injection in the input. These findings suggest that our approach is particularly effective on sparse matrix datasets.

5.1 Measuring Bias on Observed set

We begin by comparing the bias between \widehat{T} and \overline{T} , measured by the sum of squared errors between \widehat{T} (and \overline{T}) with T on Ω . Our finding is illustrated in Figure 4. First, we find that the bias of \widehat{T} compared to T, is at the order of 10^{-4} to 10^{-2} for various p. In Figure 4a, we find that on synthetic datasets, the bias of \hat{T} is 9×10^{-4} and the bias of \overline{T} is 0.3 averaged over various sampling probabilities, resulting in a reduction of 99%. In Figure 4b, we observe that for three MovieLens datasets, the bias of \hat{T} is 6×10^{-4} and the bias of \overline{T} is 5×10^{-3} , resulting in a reduction of 88%.

Additionally, we test a biased sampling procedure, where users who watch more movies tend to rate more movies as well (also known as "snowball effects" (Chen et al., 2020)). For each $i = 1, \ldots, n$, we sample its entries with probability p_i equal to C/d times the number of non-zeros in row i. We find that under this biased sampling procedure, the bias of \hat{T} is 3×10^{-4} and the bias of \overline{T} is 5×10^{-3} , resulting in a reduction of 93%. As explained earlier, these results stem from the variance reduction effect of \hat{T} .



Figure 5: Illustration of the recovery error of Algorithm 1. 5a: We vary the sampling probability and find that the error roughly stays constant after we also fix the number of samples. 5b: We see similar results as we fix the number of samples but vary the rank of the underlying matrix. 5c: We gradually increase the noise level and find that the estimation errors also increase. The reported results are based on five runs.

5.2 Recovery on Synthetic Data

We now examine one-sided recovery error by keeping the number of samples m fixed while varying p and d separately. Since the sample complexity on n grows linearly with d times $\log(d)$, the error should remain roughly constant as we fix m/d while slowly increasing d. We report simulation results under three regimes. For all three setting, we fix the number of rows $n = 10^4$.

First, we set the rank r = 10, noise injection variance $\sigma^2 = 1$, and vary d from 10^3 to 5×10^3 and p from 2/d to 10/d. As shown in Figure 5a, the estimation error remains flat for different values of p and m, which suggests that sample complexity does not grow with d, after adjusting for m/d.

Second, we set the sampling probability p = 2/d, $\sigma^2 = 1$, and vary d from 10^3 to 5×10^3 and r from 5 to 25. As shown in Figure 5b, the estimation error remains flat for a given rank r, as d grows.

Third, we set the dimension $d = 10^3$, p = 2/d, r = 10, while varying σ^2 from 1 to 50. We find that as σ^2 increases, the estimation error grows linearly, as shown in Figure 5c. This is consistent with our bounds in Theorem 3.4, where n grows proportionally with σ^2 .

Similar results on synthetic data are further illustrated in Figure 6 (See Appendix C.1). We see a consistent trend in which our approach reduces the recovery error across different settings of C (expected number of entries at each row) and n, compared to all three baseline methods, including alternating-GD, softImpute-ALS, and nuclear norm regularization. In the setting with the fewest observed entries (C = 2 and $n = 10^4$), HÁJEK-GD achieves the lowest error of 0.10, representing a 70% reduction compared to the closest baseline, nuclear norm regularization, which yields an error of 0.48. In the setting with C = 10 observed entries per row, HÁJEK-GD achieves an error of 0.06, outperforming alternative-GD and nuclear norm regularization, whose error rates are 0.09 and 0.17, respectively.

5.3 Recovery on Real-World Data

Next, we report the results on three real-world datasets in Table 2. We find that our approach consistently achieves the lowest error for recovering T and M, compared to two baseline methods. For one-sided matrix completion, HÁJEK-GD improves upon alternating-GD by 42% and softImpute-ALS by 59% on average. For recovering the entire matrix, our approach reduces recovery error by 21% relative to alternating-GD and by 38% relative to softImpute-ALS, averaged over the three datasets. In particular, on the Amazon reviews dataset, HÁJEK-GD reduces estimation of T by 59% and recovery of M by 21% relative to the best performing baseline.

An important design in HÁJEK-GD is the use of an incoherence regularization penalty of $R(\cdot)$ (cf. Section 3.2). We now report ablation studies of varying the regularization parameter λ and the threshold α for

Table 2: Results from applying our approach on three real-world datasets with up to 32 million entries, as compared with two baseline methods, namely alternating-GD and softImpute-ALS. We report both the recovery errors and the running time (in seconds, evaluated on an Ubuntu server). We run each experiment with five random seeds to calculate the mean and the standard deviation.

Dataset	MovieLens-32M	Amazon Reviews	Genomes	MovieLens-32M	Amazon Reviews	Genomes
T	One-sided recovery error			Running time		
AlternatingGD SoftImputeALS Algorithm 1	$ \begin{vmatrix} 9.1_{\pm 0.1} \times 10^{-3} \\ 9.1_{\pm 0.1} \times 10^{-3} \\ 4.7_{\pm 0.1} \times 10^{-3} \end{vmatrix} $	$\begin{array}{c} 3.3_{\pm 1.3} \times 10^{-3} \\ 3.2_{\pm 0.3} \times 10^{-3} \\ 1.3_{\pm 0.2} \times 10^{-3} \end{array}$	$\begin{array}{c} 7.0_{\pm 1.6}\times 10^{-5}\\ 1.9_{\pm 0.5}\times 10^{-4}\\ \textbf{5.8}_{\pm 0.1}\times 10^{-5}\end{array}$	$\begin{array}{c c} 7.2_{\pm 0.1} \times 10^2 \\ 6.3_{\pm 4} \times 10^4 \\ 1.5_{\pm 0.1} \times 10^2 \end{array}$	$\begin{array}{c} 4.8_{\pm 0.1} \times 10^2 \\ 5.1_{\pm 0.3} \times 10^3 \\ \textbf{2.8}_{\pm 0.2} \times 10^2 \end{array}$	$\begin{array}{c} 23.2_{\pm 5.0} \\ 959.8_{\pm 40} \\ 1.3_{\pm 0.1} \end{array}$
M	Root mean squared recovery error		Running time			
AlternatingGD SoftImputeALS Algorithm 2	$\begin{array}{c c} 1.7_{\pm 0.1} \\ 1.4_{\pm 0.1} \\ 1.1_{\pm 0.1} \end{array}$	$\begin{array}{c} 2.6_{\pm 0.1} \\ 3.3_{\pm 0.1} \\ 1.9_{\pm 0.1} \end{array}$	$\begin{array}{c} 0.2_{\pm 0.1} \\ 0.4_{\pm 0.1} \\ 0.2_{\pm 0.1} \end{array}$	$ \begin{array}{c c} 5.7_{\pm 0.1} \times 10^2 \\ 6.3_{\pm 4} \times 10^4 \\ \textbf{5.6}_{\pm 0.1} \times 10^2 \end{array} $	$\begin{array}{c} 2.2_{\pm 0.1} \times 10^2 \\ 4.8_{\pm 0.3} \times 10^3 \\ \textbf{2.8}_{\pm 0.2} \times 10^2 \end{array}$	$\begin{array}{c} 23.2_{\pm 5.0} \\ 959.9_{\pm 40} \\ 14.7_{\pm 0.2} \end{array}$

a synthetic dataset with p = 10/d and also the Amazon Review dataset. We consider the regularization parameter λ between 10^{-4} , 10^{-3} , 10^{-2} , and α (used in the truncation of $R(\cdot)$) between 10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} . We find that on synthetic datasets, the lowest recovery error is achieved when $\lambda = 10^{-4}$ and $\alpha = 10^{-3}$ at 2.8×10^{-3} , compared to 3.6×10^{-3} without using the regularizer. In other words, we can reduce the recovery error by 22%.

On the Amazon Reviews dataset, the lowest recovery error is achieved when $\lambda = 10^{-2}$ and $\alpha = 10^{-1}$, at 1.3×10^{-2} , compared to 2.7×10^{-2} without using the regularizer. As a result, incorporating incoherence regularization reduces error by 52%.

6 Related Work

The problem of recovering a low-rank matrix from a few potentially noisy entries has a rich history of studies in the literature, spanning a wide variety of areas, including collaborative filtering, machine learning, and compressed sensing (Candes & Plan, 2010). One of the earliest results from the matrix completion literature involves a maximum-margin matrix factorization approach (Srebro et al., 2004), which minimizes the reconstruction error plus a trace norm penalty on the low-rank factors. Under an incoherent condition on the low-rank factors plus a joint incoherence condition, Candes & Recht (2012) show that exact recovery of the unknown matrix is possible by minimizing the nuclear norm of the reconstructed matrix, subject to equality constraints on the observed entries. The joint incoherence condition on the factors is further shown to be unnecessary (Chen, 2015). Ding & Chen (2020) sharpen the known bounds through a leave-one-out analysis, which leads to a convergence guarantee for projected gradient descent on a rank-constrained formulation. Complementary to these results, Zhang et al. (2019b) show error bounds on the nuclear norm regularized objective, relative to the error of the best rank-r approximation of M.

Recent literature has studied the optimization geometry of the loss landscape of matrix completion using rank-constrained first-order optimization methods (Sun & Luo, 2016; Ge et al., 2016). Provided with a regularization on the incoherence (balance between the rows of low-rank factors), it is possible to prove that all local minimum solutions of the matrix completion reconstruction error objective are also global minimum (Ge et al., 2017). By carefully analyzing the local geometry of the matrix factorization objective, it is possible to achieve exact recovery under the incoherence condition (Sun & Luo, 2016). Provided with this characterization, one can obtain convergence rates through the literature on finding local minima in nonconvex optimization problems (Nesterov & Polyak, 2006; Wang et al., 2019). In particular, it is known that many optimization algorithms, including cubic regularization, trust-region methods, and stochastic gradient descent, can efficiently find a local minimizer. See a recent textbook on first-order and stochastic optimization methods for further references (Lan, 2020).

There is a closely related problem of low-rank matrix recovery, whose aim is to reconstruct a low-rank matrix based on a linear system of measurement equations of the unknown matrix (Recht et al., 2010). The

optimization landscape of low-rank matrix recovery from random linear measurements in the presence of arbitrary outliers is studied in Li et al. (2020). Besides, there are studies on the dynamics of gradient descent assuming the factorization is over-parameterized (Li et al., 2018; Ma & Fattahi, 2024). Additionally, matrix factorization has connections to two-layer neural networks and random matrices, which have inspired studies on the dynamics of deep linear networks (Pennington & Worah, 2017; Hu et al., 2020). Another related problem is nonnegative matrix factorization, which has also been studied for "tall-and-skinny" matrices (Benson et al., 2014). There is also a line of work on the recovery of a low-rank plus sparse matrix (Hsu et al., 2011). Another plausible direction in the ultra-sparse sampling setting is to instead recover a low-rank approximation of the underlying matrix from just O(n) samples. See Gamarnik et al. (2017) for further references related to this setting.

The one-sided matrix completion problem has been recently formulated and studied in the special case where each row contains only two randomly observed entries (Cao et al., 2023). This work complements their paper in four aspects. First, their sample complexity bound focuses on the setting where every row has two observed entries, while our result applies to more general settings that allow for any constant C. Second, their result builds on proof techniques for matrix completion with nuclear norm regularization penalty. By contrast, our result applies to gradient-based optimization algorithms, which are faster and easier to implement in practice. Third, our result builds on the incoherence assumption from the MC literature. Finally, we explore the use of one-sided matrix completion for full matrix completion and find that this can also lead to improved results for imputing missing entries of the full matrix.

First-order stochastic gradient methods also have applications in large-scale matrix completion (Mackey et al., 2011). In particular, one can apply stochastic gradient updates in an asynchronous protocol, making it suitable for distributed platforms (Recht & Ré, 2013). Besides the nuclear norm minimization approach, alternating minimization such as alternating least squares and low-rank matrix factorization are also widely used in practice (Hastie et al., 2015). More recently, Wang et al. (2023) consider a low-rank matrix factorization approach to private matrix completion. Assuming the M matrix is indeed of low rank, their work shows a sample complexity bound that scales linearly in dimension. The difference between our work and this work is that we focus on an ultra-sparse sampling regime where there are only O(n) randomly-sampled entries, rendering the accurate recovery of the entire matrix information-theoretically not possible. In light of our new results, it may also be worth resisting the dynamics of zeroth-order and first-order methods for nonconvex optimization (Ghadimi & Lan, 2013; Sun, 2020) when high rates of noise are added during each iteration.

The idea of using the estimated probability to reweight a population statistic is known as Hájek's estimator (Hájek, 1971). It has often been used to tackle sparse and non-random sampling data (Särndal et al., 2003). It is known from the causal inference literature (Hirano et al., 2003) that Hájek's estimator incurs lower variance than the Horvitz-Thompson estimator, which weights each observation inversely with the true probability. To our knowledge, the analysis of this estimator for one-sided matrix completion appears to be new in the MC literature. Unlike conventional results (Hájek, 1971) where the Hájek estimator is shown to be asymptotically unbiased, for one-sided matrix completion, we find that it is unbiased even in the finite sample regime. Recent work (Bai & Ng, 2021) develop the estimation of counterfactuals when potential outcomes have a factor structure for block-missing panel data. Xiong & Pelger (2023) and Duan et al. (2023) develop the inferential theory for latent factor models in large-dimensional panel data with more general non-random missing patterns by proposing a PCA-based estimator on an adjusted covariance matrix.

7 Conclusion

In this paper, we study the problem of matrix completion in an ultra-sparse sampling regime, where in each row, only a constant number of entries are observed (on average). While recovering the entire matrix is not possible, we focus on the one-sided matrix completion problem. We apply Hájek's estimator from the econometrics literature to this problem, and present a thorough analysis of this estimator. We demonstrate that this estimator is unbiased and achieves a lower variance compared to the Horvitz-Thompson estimator. Then, we use gradient descent to impute the missing entries of the second-moment matrix and analyze its sample complexity under a low-rank mixture model. The sample complexity bound is optimal in its dependence on the dimension d. Our result applies to the case of observing C entries at each row, resolving an open question from prior work.

Extensive experiments on both synthetic and real-world datasets demonstrate that our approach is more efficient for sparse matrix datasets compared to several commonly used matrix completion methods. Ablation studies validate the use of a low-rank matrix incoherence regularizer in the algorithm.

Our work opens up several avenues for future work. First, it would be interesting to better understand the sensitivity of Hájek's estimator pertaining to white noise injection, which may facilitate the design of private matrix completion. Second, the sample complexity bound has a high dependence on rank. It would be interesting to see if this dependence can be improved. Finally, it would also be interesting to revisit tensor completion in an ultrasparse setting.

References

- Susan Athey, Mohsen Bayati, Nikolay Doudchenko, Guido Imbens, and Khashayar Khosravi. Matrix completion methods for causal panel data models. *Journal of the American Statistical Association*, 116 (536):1716–1730, 2021. 8, 32
- Jushan Bai and Serena Ng. Matrix completion, counterfactuals, and factor analysis of missing data. *Journal* of the American Statistical Association, 116(536):1746–1763, 2021. 14
- Austin R Benson, Jason D Lee, Bartek Rajwa, and David F Gleich. Scalable methods for nonnegative matrix factorizations of near-separable tall-and-skinny matrices. Advances in Neural Information Processing Systems, 2014. 14
- Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. Communications of the ACM, 55(6):111–119, 2012. 1, 2, 3, 13
- Emmanuel J Candes and Yaniv Plan. Matrix completion with noise. Proceedings of the IEEE, 98(6):925–936, 2010. 1, 13
- Steven Cao, Percy Liang, and Gregory Valiant. One-sided matrix completion from two observations per row. International Conference on Machine Learning, 2023. 2, 3, 7, 10, 14, 34
- Justin Chen, Gregory Valiant, and Paul Valiant. Worst-case analysis for randomly collected data. Advances in Neural Information Processing Systems, 33:18183–18193, 2020. 11
- Yudong Chen. Incoherence-optimal matrix completion. *IEEE Transactions on Information Theory*, 61(5): 2909–2923, 2015. 2, 13
- Yuejie Chi, Yue M Lu, and Yuxin Chen. Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2019. 1, 3, 11
- Lijun Ding and Yudong Chen. Leave-one-out approach for matrix completion: Primal and dual analysis. IEEE Transactions on Information Theory, 66(11):7274–7301, 2020. 13
- Edgar Dobriban and Yue Sheng. Wonder: Weighted one-shot distributed ridge regression in high dimensions. Journal of Machine Learning Research, 21(66):1–52, 2020.
- Junting Duan, Markus Pelger, and Ruoxuan Xiong. Target pca: Transfer learning large dimensional panel data. Journal of Econometrics, pp. 105521, 2023. 14
- Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science, 9(3-4):211-407, 2014. 10
- David Gamarnik, Quan Li, and Hongyi Zhang. Matrix completion from O(n) samples in linear time. In Conference on Learning Theory, 2017. 3, 14
- Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. Advances in Neural Information Processing Systems, 2016. 3, 7, 9, 13, 24, 31

- Rong Ge, Chi Jin, and Yi Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, 2017. 7, 13, 28
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. SIAM Journal on Optimization, 23(4):2341–2368, 2013. 14
- Jaroslav Hájek. Comment on "an essay on the logical foundations of survey sampling, part one". The foundations of survey sampling, 236, 1971. 2, 4, 14
- Moritz Hardt and Aaron Roth. Beating randomized response on incoherent matrices. In Proceedings of the forty-fourth annual ACM symposium on Theory of computing, pp. 1255–1268, 2012. 2
- Trevor Hastie, Rahul Mazumder, Jason D Lee, and Reza Zadeh. Matrix completion and low-rank svd via fast alternating least squares. *The Journal of Machine Learning Research*, 16(1):3367–3402, 2015. 3, 10, 14, 33
- Keisuke Hirano, Guido W Imbens, and Geert Ridder. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71(4):1161–1189, 2003. 2, 4, 14
- Daniel G Horvitz and Donovan J Thompson. A generalization of sampling without replacement from a finite universe. Journal of the American statistical Association, 47(260):663–685, 1952. 2, 4
- Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi Chen, and Julian McAuley. Bridging language and items for retrieval and recommendation. arXiv preprint arXiv:2403.03952, 2024. 2, 10
- Daniel Hsu, Sham M Kakade, and Tong Zhang. Robust matrix decomposition with sparse corruptions. IEEE Transactions on Information Theory, 57(11):7221–7234, 2011. 14
- Wei Hu, Lechao Xiao, and Jeffrey Pennington. Provable benefit of orthogonal initialization in optimizing deep linear networks. *International Conference on Learning Representations*, 2020. 14
- Mladen Kolar, John Lafferty, and Larry Wasserman. Union support recovery in multi-task learning. Journal of Machine Learning Research, 12(7), 2011. 7
- Guanghui Lan. First-order and stochastic optimization methods for machine learning, volume 1. Springer, 2020. 2, 13
- Yuanxin Li, Yuejie Chi, Huishuai Zhang, and Yingbin Liang. Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent. Information and Inference: A Journal of the IMA, 9(2):289–325, 2020. 3, 14
- Yuanzhi Li, Tengyu Ma, and Hongyang R. Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, 2018. 14
- Ziqi Liu, Yu-Xiang Wang, and Alexander Smola. Fast differentially private matrix factorization. In Proceedings of the 9th ACM Conference on Recommender Systems, pp. 171–178, 2015. 2, 9
- Jianhao Ma and Salar Fattahi. Convergence of gradient descent with small initialization for unregularized matrix completion. In Conference on Learning Theory, 2024. 14
- Lester Mackey, Michael Jordan, and Ameet Talwalkar. Divide-and-conquer matrix factorization. Advances in Neural Information Processing Systems, 2011. 14
- Yurii Nesterov and Boris T Polyak. Cubic regularization of newton method and its global performance. Mathematical Programming, 108(1):177–205, 2006. 13
- Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. Advances in Neural Information Processing Systems, 2017. 14
- Benjamin Recht. A simpler approach to matrix completion. Journal of Machine Learning Research, 12(12), 2011. 7

- Benjamin Recht and Christopher Ré. Parallel stochastic gradient algorithms for large-scale matrix completion. Mathematical Programming Computation, 5(2):201–226, 2013. 14
- Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010. 13
- Carl-Erik Särndal, Bengt Swensson, and Jan Wretman. Model assisted survey sampling. Springer Science & Business Media, 2003. 2, 8, 14
- Nathan Srebro, Jason Rennie, and Tommi Jaakkola. Maximum-margin matrix factorization. Advances in Neural Information Processing Systems, 2004. 1, 13
- Ruo-Yu Sun. Optimization for deep learning: An overview. Journal of the Operations Research Society of China, 8(2):249–294, 2020. 14
- Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016. 2, 3, 13
- Lingxiao Wang, Boxin Zhao, and Mladen Kolar. Differentially private matrix completion through low-rank matrix factorization. In *International Conference on Artificial Intelligence and Statistics*, 2023. 2, 9, 14
- Zhe Wang, Yi Zhou, Yingbin Liang, and Guanghui Lan. Stochastic variance-reduced cubic regularization for nonconvex optimization. In The 22nd International Conference on Artificial Intelligence and Statistics, 2019. 13
- Ruoxuan Xiong and Markus Pelger. Large dimensional latent factor modeling with missing observations and applications to causal inference. *Journal of Econometrics*, 233(1):271–301, 2023. 14, 32
- Hongyang R. Zhang, Vatsal Sharan, Moses Charikar, and Yingyu Liang. Recovery guarantees for quadratic tensors with limited observations. In *International Conference on Artificial Intelligence and Statistics*, 2019a. 34
- Lijun Zhang, Tianbao Yang, Rong Jin, and Zhi-Hua Zhou. Relative error bound analysis for nuclear norm regularized matrix completion. *Journal of Machine Learning Research*, 20(97):1–22, 2019b. 10, 13

A Proof of Theorem 3.2

First, we present the proof of the unbiasedness of Hájek's estimator on the observed set Ω .

Proof of Lemma 3.1. Suppose that $(I^{\top}I)_{i,j} = k$, for any $k \ge 1$. Then, conditioned on $(I^{\top}I)_{i,j} = k$, the expectation of $\widehat{T}_{i,j}$ is equal to

$$\mathbb{E}\left[\widehat{T}_{i,j}\Big|(I^{\top}I)_{i,j}=k\right] = \mathbb{E}\left[\frac{1}{k}(\widehat{M}^{\top}\widehat{M})_{i,j}\Big|(I^{\top}I)_{i,j}=k\right]$$
$$= \mathbb{E}\left[\frac{1}{k}\sum_{a=1}^{n}M_{a,i}M_{a,j}I_{a,i}I_{q,j}\Big|(I^{\top}I)_{i,j}=k\right]$$
$$= \frac{1}{k}\sum_{a=1}^{n}\Pr\left[I_{a,i}I_{a,j}=1\Big|(I^{\top}I)_{i,j}=k\right]\cdot M_{a,i}M_{a,j}.$$

Notice that for a fixed a, the event that $I_{a,i}I_{a,j} = 1$ are pairwise independent for all a = 1, 2, ..., n. Given that k out of the n pairs are chosen, the probability that q is chosen is equal to $\binom{n-1}{k-1} / \binom{n}{k} = k/n$. Thus, the above is equal to

$$\frac{1}{k}\sum_{a=1}^{n}\frac{k}{n}\cdot M_{a,i}M_{a,j}=T_{i,j},$$

which concludes the proof of equation (3). To finish the proof that \hat{T} is unbiased, notice that

$$\mathbb{E}\left[\widehat{T}_{i,j}\Big| (I^{\top}I)_{i,j} \neq 0\right] = \sum_{k=1}^{d} \Pr\left[(I^{\top}I)_{i,j} = k \Big| (I^{\top}I)_{i,j} \neq 0 \right] \cdot \mathbb{E}\left[\widehat{T}_{i,j} \Big| (I^{\top}I)_{i,j} = k \right]$$
$$= \sum_{k=1}^{d} \Pr\left[(I^{\top}I)_{i,j} = k \Big| (I^{\top}I)_{i,j} \neq 0 \right] \cdot T_{i,j} = T_{i,j},$$

which concludes the proof of equation 2.

Next, we present the proof of Theorem 3.2, which approximates the variance of Hjake's estimator on Ω .

Proof of Theorem 3.2. For equation (4), we use Taylor's expansion on $Var(\widehat{T}_{i,i})$. For ease of presentation, let us denote $A_{i,i} = (\widehat{M}^{\top} \widehat{M})_{i,i}$, and let $B_{i,i} = (I^{\top} I)_{i,i}$, for any $1 \leq i \leq d$. Thus, conditioned on the event that $B_{i,i} \neq 0$, the (i, i)-th entry of \widehat{T} , given by $\widehat{T}_{i,i}$, is equal to $\frac{A_{i,i}}{B_{i,i}}$. By Taylor's expansion, we approximate the ratio statistic using the delta method as:

$$\frac{A_{i,i}}{B_{i,i}} = \frac{\mathbb{E}\left[A_{i,i}\right]}{\mathbb{E}\left[B_{i,i}\right]} + \frac{1}{\mathbb{E}\left[B_{i,i}\right]} \left(A_{i,i} - \mathbb{E}\left[A_{i,i}\right]\right) - \frac{\mathbb{E}\left[A_{i,i}\right]}{(\mathbb{E}\left[B_{i,i}\right])^2} \left(B_{i,i} - \mathbb{E}\left[B_{i,i}\right]\right) + \epsilon_{i,i},$$
(18)

where $\epsilon_{i,i}$ is at the order of the variance of $A_{i,i}$ and $B_{i,i}$, which is O(1/(np)). Thus, based on the first-order expansion, the variance of $\frac{A_{i,i}}{B_{i,i}}$ is approximated by:

$$Var\left(\frac{A_{i,i}}{B_{i,i}}\right) = \frac{Var\left(A_{i,i}\right)}{(\mathbb{E}\left[B_{i,i}\right])^2} + \frac{(\mathbb{E}\left[A_{i,i}\right])^2 Var\left(B_{i,i}\right)}{(\mathbb{E}\left[B_{i,i}\right])^4} - \frac{2\mathbb{E}\left[A_{i,i}\right] \cdot Cov(A_{i,i}, B_{i,i})}{(\mathbb{E}\left[B_{i,i}\right])^3} + O\left(\frac{1}{np}\right),$$
(19)

where the error term is based on the approximation error that stems from $\epsilon_{i,i}$ in equation (18).

Now we look into simplifying equation (19). Based on the definition of $A_{i,i}, B_{i,i}$, we get the following facts, for any $i = 1, \ldots, d$:

$$\mathbb{E}[B_{i,i}] = np, \quad Var(B_{i,i}) = np(1-p), \tag{20}$$

$$\mathbb{E}[A_{i,i}] = p(M^{\top}M)_{i,i}, \quad Cov(A_{i,i}, B_{i,i}) = (p - p^2)(M^{\top}M)_{i,i}.$$
(21)

To verify the above results, let us introduce a Bernoulli random variable X_k , which is equal to 1 with probability p, or 0 with probability 1-p, for k = 1, 2, ..., n. Thus, we can write

$$A_{i,i} = \sum_{k=1}^{n} M_{k,i}^2 X_k$$
, and $B_{i,i} = \sum_{k=1}^{n} X_k$.

As a result,

$$\begin{aligned} Cov(A_{i,i}, B_{i,i}) &= \mathbb{E} \left[A_{i,i} B_{i,i} \right] - \mathbb{E} \left[A_{i,i} \right] \cdot \mathbb{E} \left[B_{i,i} \right] \\ &= \mathbb{E} \left[\left(\sum_{k=1}^{n} M_{k,i}^{2} X_{k} \right) \left(\sum_{k=1}^{n} X_{k} \right) \right] - np^{2} (M^{\top} M)_{i,i} \\ &= \mathbb{E} \left[\left(\sum_{k=1}^{n} M_{k,i}^{2} X_{k} \right) + \left(\sum_{1 \le k \ne k' \le n} M_{k,i}^{2} X_{k} X_{k'} \right) \right] - np^{2} (M^{\top} M)_{i,i} \quad \text{(note that } X_{k}^{2} = X_{k}) \\ &= \left(\sum_{k=1}^{n} M_{k,i}^{2} \right) (p + (n-1)p^{2}) - np^{2} (M^{\top} M)_{i,i} \\ &= (p - p^{2}) (M^{\top} M)_{i,i}. \end{aligned}$$

Therefore, by plugging in the results from equations (20) and (21) back into equation (19), we obtain an approximate variance of the ratio statistic as:

$$\frac{Var(A_{i,i})}{n^2 p^2} + \frac{p^2 (M^\top M)_{i,i}^2 \cdot np(1-p)}{n^4 p^4} - \frac{2p (M^\top M)_{i,i} \cdot (p-p^2) (M^\top M)_{i,i}}{n^3 p^3} \\
= \frac{Var(A_{i,i})}{n^2 p^4} - \frac{(1-p) (M^\top M)_{i,i}^2}{n^3 p}.$$
(22)

This applies to any diagonal entries of \widehat{T} . In particular, notice that the first term of equation (22) is precisely the variance of $\overline{T}_{i,i}$, which reweights the observed entries based on the true probability. This leads to the following variance estimate for the diagonal entries:

$$\frac{Var(A_{i,i})}{n^2 p^2} - \frac{(1-p)(M^{\top}M)_{i,i}^2}{n^3 p}.$$
(23)

Finally, we write down the variance of $A_{i,i}$ as follows:

$$Var(A_{i,i}) = p(1-p) \sum_{k=1}^{n} M_{k,i}^{4}.$$
(24)

Notice that the probability that $B_{i,i}$ is equal to zero is at most $O(d^{-C})$. We can add this error analysis into the above proof, which does not affect the order of the error term. This concludes the proof of equation (4) in Theorem 3.2 regarding the variance approximation of the diagonal entries.

Next, we consider the case of off-diagonal entries in equation (5). The probability that an off-diagonal entry of \hat{T} is nonzero is

$$1 - (1 - p^2)^n$$
,

which is approximately $np^2 = C^2(\log d)^2/d$. Thus, conditioned on the fact that $(I^{\top}I)_{i,j} \neq 0$, the dominating event is when $(I^{\top}I)_{i,j} = 1$. The error to this, which is when there are at least two nonzero entries in the sum, is less than

$$\frac{\Pr[(I^{\top}I)_{i,j} \ge 2]}{\Pr[(I^{\top}I)_{i,j} \ne 0]} \approx np^2$$

When the (i, j)-th entry of $I^{\top}I$ is equal to one, then $\widehat{T}_{i,j} = (\widehat{M}^{\top}\widehat{M})_{i,j}$. The variance of $\widehat{T}_{i,j}$, conditioned on $(I^{\top}I)_{i,j} = 1$, is thus given by

$$\frac{1}{n}\sum_{k=1}^{n} \left(M_{k,i}M_{k,j} - \frac{1}{n} \left(\sum_{k'=1}^{n} M_{k',i}M_{k',j} \right) \right)^{2}.$$

After simplifying the above, we reach the conclusion of equation (5). The proof is thus completed.

B Proof of Theorem 3.4

We analyze the sample complexity of Algorithm 1. In particular, we will analyze $\ell(\cdot)$. Without loss of generality, let us assume that $||U||_F^2 = r$ within this section. This also implies $\sigma_{\max}(U) \ge 1 \ge \sigma_{\min}(U)$.

A key part of the proof is analyzing the bias of \hat{T} (relative to T). Let $N \in \mathbb{R}^{d \times d}$ denote the bias of Hájek's estimator, where

$$N_{i,j} = T_{i,j} - T_{i,j}, \forall (i,j) \in \Omega; \text{otherwise}, N_{i,j} = 0, \forall (i,j) \notin \Omega.$$

Here is a road map of this section.

• First, in Appendix B.1, we bound the variance of Hájek's estimator, on the observed set Ω .

- Second, in Appendix B.2, we derive a concentration bound that applies to the weighted projection P_{Ω} . By a similar argument, we also derive the spectral norm of N.
- Then, in Appendix B.3, we derived the local optimality conditions from $\ell(\cdot)$, and then we argue that any local minimum solution of $\ell(\cdot)$ must satisfy the incoherence assumption with a suitable condition number.
- Next, in Appendix B.4, we use the incoherence assumption to reduce the local optimality conditions on the empirical loss to the population loss.
- Finally, in Appendix B.5, we analyze the population loss and based on that, complete the proof of Theorem 3.4.

B.1 Proof of Lemma 4.1

In the following proof, we bound the variance of the diagonal entries of \hat{T} and also its off-diagonal entries grouped by each row.

Proof. Based on Lemma 3.1, we already know that \hat{T} is unbiased on Ω . Next, we derive the variance of the diagonal and off-diagonal entries of \hat{T} separately, since

$$\mathbb{E}\left[\left(\widehat{T}_{i,j} - T_{i,j}\right)^2\right] = Var(\widehat{T}_{i,j}), \text{ for any } (i,j) \in \Omega.$$
(25)

For the diagonal entries of T, with probability at least $1 - O(d^{-1})$, all of them are observed in Ω . In particular, when $n \ge d \log d/\epsilon^2$, by the Chernoff bound and union bound, with probability at least $1 - O(d^{-1})$, for all $i = 1, 2, \ldots, d$,

$$\left| (I^{\top}I)_{i,i} - \mathbb{E}\left[(I^{\top}I)_{i,i} \right] \right| \le \epsilon \sqrt{\frac{6C^{-1}}{d}} np.$$

Thus, we can get that $Var(\hat{T}_{i,i}) \leq (1 + 3\epsilon \sqrt{6C^{-1}/d}) Var(\frac{A_{i,i}}{np})$. Next, notice that

$$A_{i,i} = \sum_{k=1}^{n} M_{k,i}^2 X_k.$$
 (26)

Therefore, the variance of $A_{i,i}/(np)$ is equal to

$$\begin{aligned} Var\left(\frac{A_{i,i}}{np}\right) &= \frac{1-p}{n^2 p} \sum_{k=1}^n M_{k,i}^4 \\ &\leq \frac{1-p}{np} \left(\frac{1}{r} \sum_{s=1}^r u_{i,s}^4\right) \left(1 + \frac{\max_{1 \leq s \leq r} u_{i,s}^4}{r^{-1} \sum_{s=1}^r u_{i,s}^4} \sqrt{\frac{3\log(d/2)}{2n}}\right) \end{aligned}$$
(by Hoeffding's inequality)
$$&\leq \frac{1-p}{np} \left(\frac{1}{r} \left(\sum_{s=1}^r u_{i,s}^2\right)^2\right) \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}}\right) \\ &\leq \frac{1-p}{np} \frac{\mu^2 r}{d^2} \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}}\right),\end{aligned}$$

which concludes the proof of equation (11). Similar to this calculation, by applying Hoeffding's inequality, we can also show that

$$\frac{A_{i,i}}{np} \le \frac{1}{r} \sum_{s=1}^{r} \mu_{i,s}^2 \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}} \right) \le \frac{\mu}{d} \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}} \right).$$
(27)

Next, we examine the off-diagonal entries of \hat{T} . From Theorem 3.2, the variance of the off-diagonal entries is approximated by the variance of a single entry. Consider the case where we condition on taking two overlapping entries from columns *i* and *j*. Denote the two entries as a_1, a_2 . We have that

$$Var\left(\frac{a_1+a_2}{2}\right) = \frac{Var(a_1)+Var(a_2)}{4} + Cov(a_1,a_2) \le \frac{Var(a_1)+Var(a_2)}{2},$$

by the Cauchy-Schwarz inequality on the covariance. The probability that we observe three overlapping entries is less than $\tilde{O}(d^{-3})$. Since the size of Ω is $O(d^2q) = \tilde{O}(d^{-1})$. The probability that such events can happen is less than $\tilde{O}(d^{-1})$. In summary, we have shown that for every $(i, j) \in \Omega$ where $i \neq j$,

$$Var(\widehat{T}_{i,j}) \le \left(\frac{1}{n} \left(\sum_{k=1}^{n} M_{k,i}^2 M_{k,j}^2\right) - T_{i,j}^2\right) + \tilde{O}(d^{-3}) \le \frac{1}{n} \sum_{k=1}^{n} M_{k,i}^2 M_{k,j}^2 + \tilde{O}(d^{-3})$$

We now look at the sum of the variances of all the observed entries in the *i*-th row. Let the set be denoted by $S_i = \{j : (i, j) \in \Omega, j \neq i\}$, for all i = 1, 2, ..., d. We have

$$\frac{1}{n} \sum_{j \in S_i} \sum_{k=1}^n M_{k,i}^2 M_{k,j}^2 \tag{28}$$

$$\leq \sum_{j \in S_i} \left(\frac{1}{r} \sum_{s=1}^r u_{i,s}^2 u_{j,s}^2 \right) \left(1 + \frac{\max_{s=1}^r u_{i,s}^2 u_{j,s}^2}{r^{-1} \sum_{s=1}^r u_{i,s}^2 u_{j,s}^2} \sqrt{\frac{3 \log(d/2)}{2n}} \right)$$
(by Hoeffding's inequality)

$$\leq \frac{1}{r} \sum_{s=1}^{r} u_{i,s}^{2} \cdot q \|u_{s}\|^{2} \left(1 + \frac{\max_{j=1}^{d} u_{j,s}^{2}}{r^{-1} \sum_{j=1}^{d} u_{j,s}^{2}} \sqrt{\frac{3 \log(d/2)}{2d}} \right) \left(1 + r \sqrt{\frac{3 \log(d/2)}{2n}} \right) \quad \text{(by Hoeffding's inequality)}$$

$$\leq \frac{q}{r} \sum_{s=1}^{r} u_{i,s}^{2} \cdot \left(1 + r \sqrt{\frac{3 \log(d/2)}{2d}} \right) \left(1 + r \sqrt{\frac{3 \log(d/2)}{2n}} \right) \quad \text{(rolay } \|u_{s}\|^{2} \text{ by its max over } e^{\frac{1}{2}}$$

$$\leq \frac{1}{r} \cdot \max_{1 \leq s' \leq r} \|u_{s'}\|^2 \cdot \sum_{s=1}^{r} u_{i,s}^2 \cdot \left(1 + r\sqrt{\frac{3 \log(d/2)}{2d}}\right) \left(1 + r\sqrt{\frac{3 \log(d/2)}{2d}}\right) \quad (\text{relax } \|u_s\|^2 \text{ by its max over } s)$$

$$\leq \frac{\mu q}{d} \cdot \left(1 + 3r\sqrt{\frac{3 \log(d/2)}{2d}}\right) \quad (\text{by the incoherence assumption on } U)$$

The last line uses the premise that the maximum ℓ_2 norm of a factor is at most 1, and also $n \ge d$. In particular, the above Hoeffding's inequality applies uniformly to all i with probability at least $1 - O(d^{-1})$. In summary, we have shown that

$$\sum_{j \in S_i} Var(\widehat{T}_{i,j}) \leq \frac{\mu q}{d} \left(1 + 3r \sqrt{\frac{3\log(d/2)}{2d}} \right) + \tilde{O}\left(\frac{nq}{d^3}\right),$$

where we use the fact that with probability at least $1 - d^{-2}$, $|S_i| \leq 2nq$. Therefore, we conclude that equation (12) is true.

As a remark, one corollary we can draw from the above calculation is that

$$Var(\widehat{T}_{i,j}) \le \frac{\mu^2 r^2}{d^2} \left(1 + r \sqrt{\frac{3\log(d/2)}{2n}} \right) + \widetilde{O}(d^{-3}).$$
 (29)

This can be seen by following the steps from equation (28) for a fixed pair of $(i, j) \in \Omega$.

Next, we build on Lemma 3.1 to examine the bias of \hat{T} on the observed set Ω . Another corollary from the above variance calculation is the following bound on the bias of \hat{T} . We consider the diagonal and off-diagonal entries separately.

Corollary B.1. In the setting of Lemma 4.1, suppose $n \ge dr$. Then, with probability at least $1 - O(d^{-1})$, for every i = 1, 2, ..., d, the following must hold:

$$N_{i,i} \le \sqrt{\frac{32C^{-1}\mu^2 \log(d)}{nd^2}},\tag{30}$$

$$\sum_{j \in S_i} |N_{i,j}| \le q \cdot \sqrt{2\mu + \tilde{O}(d^{-1/2})}.$$
(31)

Proof. First, we have that

$$\begin{aligned} |N_{i,i}| &= \left| \widehat{T}_{i,i} - T_{i,i} \right| = \left| \frac{A_{i,i}}{B_{i,i}} - T_{i,i} \right| = \left| \left(1 \pm \epsilon \sqrt{\frac{6C^{-1}}{d}} \right) \frac{A_{i,i}}{np} - T_{i,i} \right| \\ &\leq \left| \frac{A_{i,i}}{np} - T_{i,i} \right| + \epsilon \sqrt{\frac{6C^{-1}}{d}} \frac{A_{i,i}}{np}. \end{aligned}$$

By equation (27), the latter is at most

$$\epsilon \sqrt{\frac{6C^{-1}}{d}} \frac{\mu}{d} \left(1 + r\sqrt{\frac{3\log(d/2)}{2n}} \right) \le \sqrt{\frac{8C^{-1}\mu^2\log(d)}{d^2n}},$$

since $\epsilon \leq \sqrt{d \log d/n}$. As for the former, by equation (30), we have $Var(N_{i,i}) \leq (1 + \epsilon \sqrt{6C^{-1}/d}) Var(\frac{\widehat{A}_{i,i}}{np})$. Recall that $A_{i,i}$ is the sum of *n* Bernoulli random variables (cf. equation (26)), by Bernstein's inequality, with probability at least $1 - O(d^{-1})$, for all i = 1, 2, ..., d,

$$\begin{aligned} \left| \frac{A_{i,i}}{np} - T_{i,i} \right| &\leq \sqrt{\frac{6Var(N_{i,i}) \cdot \log(d/2)}{n}} + \frac{3\left(\max_{s=1}^{r} u_{i,s}^2 \right) \log(d/2)}{3n} \\ &\leq \sqrt{\frac{7\mu^2 r \log(d/2)}{d^2 n^2 p}} + \frac{\mu^2 r^2 \log(d/2)}{nd}, \end{aligned}$$

for large enough values of d. Specifically, we have used the bound on the variance of $\widehat{T}_{i,i}$ in Lemma 4.1, using equation (11). Notice that the second term is a lower-order term relative to the first. This shows that equation (30) holds for large enough values of d.

For equation (31), based on the proof of Lemma 4.1,

$$\sum_{j \in S_i} |N_{i,j}| \le \sqrt{|S_i| \cdot \sum_{j \in S_i} N_{i,j}^2}$$
(by Cauchy Schwarz)

$$\leq \sqrt{2dq \cdot \sum_{j \in S_i} N_{i,j}^2} \qquad (\text{since } |S_i| \leq 2dq)$$

$$\leq \sqrt{2dq \cdot \left(\frac{\mu q}{d} \left(1 + 3r\sqrt{\frac{3\log(d/2)}{2d}}\right) + \tilde{O}(d^{-3}) + \frac{2\mu^2 r^2}{d^2} \cdot \sqrt{\frac{3\log(d/2)}{2d}}\right)},\tag{32}$$

where the last line is by applying Hoeffding's inequality to the sequence $\{N_{i,j}^2 : j \in S_i\}$ (which are all independent from each other), and holds with probability at least $1 - d^{-1}$ over all possible *i*. In particular, the expectation on this sequence is based on equation (12) and the maximum is based on equation (29). From the last line above, we conclude that equation (31) is true.

B.2 Proof of Lemma 4.2

Proof. By the weighted operation in P_{Ω} and the condition that all the diagonals have been observed, we can cancel out the diagonal entries between $P_{\Omega}(W)$ and qW. In the rest of the following, we focus on the

off-diagonal entries. Let $x \in \mathbb{R}^d$ denote a unit vector. We have that

$$||P_{\Omega}(W) - qW||_2 = \max_{x:||x||=1} x^{\top} (P_{\Omega}(W) - qW)x.$$

Next, we expand the right-hand side above as:

$$x^{\top}(P_{\Omega}(W) - qW)x = \sum_{i=1}^{d} \left(\sum_{j \in S_{i}: j \neq i} W_{i,j}x_{i}x_{j} - q \sum_{1 \le j \le d: j \neq i} W_{i,j}x_{i}x_{j} \right)$$
$$\leq \sum_{i=1}^{d} \left| \sum_{\substack{j \in S_{i}: j \neq i}} W_{i,j}x_{i}x_{j} - q \sum_{\substack{1 \le j \le d: j \neq i}} W_{i,j}x_{i}x_{j} \right|.$$

We focus on a fixed i above. By Bernstein's inequality, with probability at least $1 - O(d^{-2})$, we have

$$e_{i} \leq \sqrt{4q(1-q) \sum_{1 \leq j \leq d: j \neq i} \left(W_{i,j}^{2} x_{i}^{2} x_{j}^{2}\right) \log(2d) + \|W\|_{\infty} |x_{i}| \max_{1 \leq j \leq d} |x_{j}|}$$

$$\leq \|W\|_{\infty} |x_{i}| \sqrt{4q(1-q) \log(2d)} + \|W\|_{\infty} |x_{i}| \max_{1 \leq j \leq d} |x_{j}|.$$

Above, we have plugged in the variance and the maximum values at the *i*-th row. As a result,

$$\sum_{i=1}^{d} e_i \le \frac{\nu}{\sqrt{d}} \sqrt{4q(1-q)\log(2d)} + \frac{\nu}{d} \sum_{i=1}^{d} |x_i| \max_{1 \le j \le d} |x_j|.$$
(33)

In particular, we have used the bound on $||W||_{\infty} \leq \nu/\sqrt{d}$ and the fact that $\sum_{i} |x_{i}| \leq \sqrt{d}$ via the Cauchy-Schwarz inequality.

Next, we notice that

$$\sum_{i=1}^{d} |x_i| \max_{1 \le j \le d} |x_j| \le 1.$$

To see this, notice that if we look at a particular pair x_i, x_j . Conditioned on a fixed $x_i^2 + x_j^2$ and $|x_i| \le \lambda$, $|x_j| \le \lambda$, $|x_i| + |x_j|$ are maximized when $x_i = x_j$. This implies that when $\sum_{i=1}^d |x_i|$ is maximized, all the non-zero x_i 's must be equal to each other. Suppose there are k of them, and they are all equal to a. Then, we must have $ka^2 \le 1$, $a \le \lambda$. As a result,

$$\sum_{i=1}^{d} |x_i| \cdot \max_j |x_j| = ka\lambda \le 1.$$

By rearranging equation (33), and using the fact that $dq \ge 4\log(2d)$ to relax the second term in equation (33), we have concluded the proof of equation (15).

With a similar proof, we can derive a bound on the spectral norm of N, stated as follows.

Corollary B.2. In the setting of Lemma 4.1, as long as $q \ge \frac{4 \log(2d)}{d}$, then with probability at least $1 - O(d^{-1})$, the spectral norm of $P_{\Omega}(N)$ satisfies:

$$\|P_{\Omega}(N)\|_{2} \le q\sqrt{\frac{32\mu^{2}r^{2}\log(2d)}{dq}} + O\left(q\sqrt{\frac{\log(d)}{nd^{2}}}\right).$$
(34)

Proof. Let $x \in \mathbb{R}^d$ be any unit vector. Note that N is a symmetric matrix. Therefore, we expand $x^\top Nx$ as follows:

$$x^{\top} N x = q \sum_{i=1}^{d} N_{i,i} x_i^2 + \sum_{1 \le i \le d} \sum_{j \in S_i : j \ne i} N_{i,j} x_i x_j.$$

By equation (30), Corollary B.1,

$$e_1 \le q \sqrt{\frac{32C^{-1}\mu^2 \log(d)}{nd^2}} \sum_{i=1}^d x_i^2 = q \sqrt{\frac{32C^{-1}\mu^2 \log(d)}{nd^2}}$$

As for e_2 , notice that $|N_{i,j}| = \left| \hat{T}_{i,j} - T_{i,j} \right| \le 2\mu r/d$. By following the proof of Lemma 4.2 above (cf. equation (33)), we have that

$$e_2 \le \frac{2\mu r}{\sqrt{d}}\sqrt{4q(1-q)\log(2d)} + \frac{2\mu r}{d}$$

Since $dq \ge 4\log(2d)$, the second term on the right is less than or equal to the first term on the right. Put together, we can see that equation (34) is true.

Another implication is the inner product of two matrices under the projection operator.

Proposition B.3. Suppose $||X||_{\infty} \cdot ||Y||_{\infty} \leq \frac{\nu}{d}$. With probability at least $1 - O(d^{-1})$ over the randomness of Ω , for any two matrices $X, Y \in \mathbb{R}^{d \times d}$, we have

$$|\langle P_{\Omega}(X), Y \rangle - q \langle X, Y \rangle| \le q \nu \sqrt{\frac{16 \log(2d)}{dq}}.$$
(35)

Proof. With probability at least $1 - O(d^{-1})$, the diagonal entries of $\langle P_{\Omega}(X), Y \rangle$ are canceled out with that of $q\langle X, Y \rangle$. As for the off-diagonal entries, let

$$e_{2} = \sum_{1 \le i \le d} \sum_{j \in S_{i}: j \ne i} X_{i,j} Y_{i,j} - q \sum_{1 \le i \le d} \sum_{1 \le j \le d: j \ne i} X_{i,j} Y_{i,j}.$$

Notice that for a fixed *i*, whether $j \in S_i$ or another $j' \in S_i$ are independent events. Thus, we apply Bernstein's inequality to each row, similar to the steps following equation (33), which leads to the bound on e_2 stated in equation (35).

B.3 Local Optimality and Incoherence

First, we derive the first-order and second-order optimality conditions of $\ell(X)$, accounting for the weighted projection P_{Ω} to offset the imbalance between diagonal and off-diagonal entries in Ω .

Proposition B.4 (See also Prop. 5.1, Ge et al. (2016)). Suppose X is a local optimum of $\ell(X)$ (cf. Eq. (8)). Then, the first order optimality condition is equivalent to

$$2P_{\Omega}(\widehat{T})X = 2P_{\Omega}(XX^{\top})X + \lambda \nabla R(X).$$
(36)

The second-order optimality condition is equivalent to

$$\forall V \in \mathbb{R}^d, \ \langle P_{\Omega}(VX^{\top} + XV^{\top}), XV^{\top} + VX^{\top} \rangle + \lambda \langle V, \nabla^2 R(X)V \rangle \ge 2 \langle P_{\Omega}(\widehat{T} - XX^{\top}), VV^{\top} \rangle.$$
(37)

Proof. The gradient formula of equation (36) can be seen by directly taking the gradient on $\ell(X)$, and separating the gradient from the squared loss vs the regularizer. As for the Hessian formula of equation (37), when V is sufficiently small, one has that $\nabla \ell(X + V) = [\nabla^2 \ell(X)]V$. Thus, we have that

$$\nabla \ell(X+V) = 2P_{\Omega}(VX^{\top} + XV^{\top} + VV^{\top} + XX^{\top})(X+V) - 2P_{\Omega}(\widehat{T})(X+V) + \lambda \nabla R(X+V).$$

By the second order optimality condition, one has that $V^{\top} \nabla \ell(X+V) \ge 0$, when V tends to zero. Plugging this into the above, we get

$$V^{\top} \nabla \ell(X+V)$$

$$= 2 \langle P_{\Omega}(VX^{\top} + XV^{\top}), XV \rangle + 2 \langle P_{\Omega}(XX^{\top}), VV^{\top} \rangle - 2 \langle P_{\Omega}(\widehat{T}), VV^{\top} \rangle + \lambda \nabla R(X+V) + O(||V||_{F}^{2})$$

$$= \langle P_{\Omega}(VX^{\top} + XV^{\top}), XV^{\top} + VX^{\top} \rangle + 2 \langle P_{\Omega}(XX^{\top} - \widehat{T}), VV^{\top} \rangle + \lambda \langle V, \nabla^{2}R(X)V \rangle + O(||V||_{F}^{2})$$

$$= \langle P_{\Omega}(VX^{\top} + XV^{\top}), VX^{\top} + XV^{\top} \rangle + 2 \langle P_{\Omega}(XX^{\top} - \widehat{T}), VV^{\top} \rangle + \lambda \langle V, \nabla^{2}R(X)V \rangle + O(||V||_{F}^{2}).$$

When V tends to zero, we have that equation (37) must be true for any $V \in \mathbb{R}^d$ based on the second order optimality condition.

Next, we show that the regularizer of R(X) leads the row norms of X to be bounded by at most 2α , which builds on the first-order optimality conditions above. As a result, we can use concentration tools to reduce the first- and second-order optimality conditions for their population versions. Based on that, we show that X must be bounded away from zero.

We begin by analyzing the effect of the regularizer. We show the following property, which results from adding the incoherence regularizer and can be derived based on the gradient of the regularizer.

Proposition B.5. The gradient of R(X) satisfies that for any $Y \in \mathbb{R}^{d \times r}$,

$$\langle \nabla R(X), Y \rangle = 4\lambda \sum_{i=1}^{d} \left(\|X_i\|^3 - \alpha \right)_+^3 \frac{\langle X_i, Y^\top e_i \rangle}{\|X_i\|}.$$

where $e_i \in \mathbb{R}^d$ is the *i*-th basis vector. As a consequence,

$$\langle (\nabla R(X))_i, X_i \rangle \ge 0$$
, for every $i = 1, 2, \dots, d$.

We first notice that because of the regularizer, any matrix X that satisfies the first-order optimality condition must also satisfy the incoherence condition.

Lemma B.6. Let S_i denote the set of observed entries in \hat{I} at the *i*-th row, for i = 1, 2, ..., d. Suppose $|S_i| \leq 2dq$, for any i = 1, 2, ..., d. In the setting of Theorem 3.4, for any X that satisfies the first order optimality of equation (36), we have

$$\max_{i=1}^{d} \|X_i\| \le \max\left(2\alpha, 2\sqrt{\mu(r+1)q/\lambda}\right).$$
(38)

Proof. Let $i^* = \arg \max_i ||X_i||$ be the index of the row that has the highest norm in X. Recall that $X_i \in \mathbb{R}^r$ refers to the *i*-th row vector of X. Suppose the *i*-th row of Ω consists of entries with index $[i] \times S_i$, where S_i is the set of indices in the *i*-th row that are observed in Ω . If $|X_{i^*}| \leq 2\alpha$, then equation (38) is true. Otherwise, let us assume that $|X_{i^*}| \geq 2\alpha$.

We will compare the i^* -th row of the left and right-hand sides of equation (36). First, we have the following identity based on the projection matrix P_{Ω} :

$$\left(P_{\Omega}(\widehat{T})X\right)_{i^{\star}} = \left(P_{\Omega}(UU^{\top} + N)X\right)_{i^{\star}} = \left(P_{\Omega}(UU^{\top})\right)_{i^{\star}}X + N_{i^{\star}}X,$$

where the sub-index on i^* refers to the i^* -th row of the enclosed matrix. Then, the ℓ_1 norm of $P_{\Omega}(UU^{\top})$ is at most

$$\left\| \left(P_{\Omega}(UU^{\top}) \right)_{i^{\star}} \right\|_{1} = \sum_{j \in S_{i^{\star}}} \left| \langle U_{i^{\star}}, U_{j} \rangle \right| \tag{39}$$

$$\leq \sum_{j \in S_{i^{\star}}} \|U_{i^{\star}}\| \cdot \|U_{j}\| \leq \sum_{j \in S_{i^{\star}}} \frac{\mu r}{d}$$
 (by the incoherence assumption on U)

$$\leq \frac{2\mu r dq}{d} = 2\mu r q. \qquad (by |S_{i^*}| \leq 2dq)$$

Second, we look at the ℓ_1 norm of the *i*^{*}-th row of *N*. By using Corollary B.1, we can bound the ℓ_2 norm of the left of equation (36) by

$$\left\| (P_{\Omega}(\widehat{T})X)_{i^{\star}} \right\| \leq \left(\left\| (P_{\Omega}(UU^{\top}))_{i^{\star}} \right\|_{1} + \left\| (P_{\Omega}(N))_{i^{\star}} \right\|_{1} \right) \cdot \max_{i=1}^{d} \|X_{i}\| \\ \leq \left(2\mu r + \sqrt{2\mu + \tilde{O}(d^{-1/2})} + \tilde{O}(n^{-1}) \right) q \|X_{i^{\star}}\|,$$
(40)

where we use equations (39), (30), and (31) above, along with the fact that p = C/n.

Next, we lower bound the norm of the right-hand side of equation (36). We have that

$$(P_{\Omega}(XX^{\top})X)_{i^{\star}} = \sum_{j \in S_{i^{\star}}} \langle X_{i^{\star}}, X_j \rangle X_j.$$

Thus, by taking an inner product with $X_{i^{\star}}$, we get

$$\langle (P_{\Omega}(XX^{\top})X)_{i^{\star}}, X_{i^{\star}} \rangle = \sum_{j \in S_{i^{\star}}} \langle X_{i^{\star}}, X_{j} \rangle^{2} \ge 0.$$

Using Proposition B.5, we obtain that

$$\langle (P_{\Omega}(XX^{\top})X)_{i^{\star}}, (\nabla R(X))_{i^{\star}} \rangle = 4\lambda \left(|X_{i^{\star}}| - \alpha \right)_{+}^{3} \frac{\sum_{j \in S_{i^{\star}}} \langle X_{i^{\star}}, X_{j} \rangle^{2}}{\|X_{i^{\star}}\|} \ge 0.$$
 (41)

It follows that

$$\begin{aligned} \left\| (P_{\Omega}(XX^{\top})X)_{i^{\star}} + (\lambda \nabla R(X))_{i^{\star}} \right\| &\geq \| (\lambda \nabla R(X))_{i^{\star}} \| & \text{(by equation (41))} \\ &= \frac{4\lambda (\|X_{i^{\star}}\| - \alpha)_{+}^{3}}{\|X_{i^{\star}}\|} \cdot \|X_{i^{\star}}\| & \text{(by Proposition B.5)} \\ &\lambda & 2 \end{aligned}$$

$$\geq \frac{\lambda}{2} \|X_{i^{\star}}\|^3 \qquad (\text{since } \|X_{i^{\star}}\| \geq 2\alpha)$$

Therefore, plugging in the equation above and the equation (40) into the first order optimality condition (36). We obtain equation (38). Thus, the proof is completed.

The condition that $|S_i| \leq 2dq$ for all *i* can be shown via Chernoff bounds (recall that $q = 1 - (1 - p^2)^n$). To see this, notice that conditioned on row *i*, the event of whether $j \in S_i$ (for any *j* between 1 and *d* but not equal to *i*) is independent from each other. In expectation, the size of $|S_i|$ should be equal to dq. In the setting of Theorem 3.4, dq is at least $2\log(d)$. Therefore, with probability at least $1 - O(d^{-2})$, $|S_i| \leq 2dq$, for all $i = 1, 2, \ldots, d$.

B.4 Reduction of Local Optimality Conditions

One consequence of the above result is the following population version of the first-order optimality condition. **Lemma B.7.** Suppose $\alpha \leq \sqrt{\mu(r+1)q/\lambda}$ and $\lambda \geq dq/8$. In the setting of Theorem 3.4, with high probability over the randomness of Ω , for any $X \in \mathbb{R}^{d \times r}$ that satisfies the first order optimality condition (36), we have that

$$\|X\|_F \le \sigma_{\max}(U)\sqrt{r} + \sqrt{\delta}, \text{ where } \delta = 3\sqrt{32\mu^2 r^3 \log(2d)/(dq)},$$

$$\tag{42}$$

$$\left\| UU^{\top}X - XX^{\top}X - \gamma \nabla R(X) \right\|_{F} \le 42\sigma_{\max}(U)\sqrt{\frac{\mu^{2}r^{4}\log(2d)}{dq}}.$$
(43)

Proof. By the incoherence assumption on U, we have that

$$\left\| UU^{\top} \right\|_{\infty} \leq \frac{\mu r}{d} \left\| U \right\|_{F}^{2} \leq \frac{\mu r^{3/2}}{d} \left\| UU^{\top} \right\|_{F}$$

by applying the Cauchy-Schwarz inequality. By setting $\nu = \mu r^{3/2} \|UU^{\top}\|_F$ in Lemma 4.2, from equation (15) we get

$$\left\|P_{\Omega}(UU^{\top})X - qUU^{\top}X\right\|_{F} \le \left\|P_{\Omega}(UU^{\top})X - qUU^{\top}\right\|_{2} \cdot \left\|X\right\|_{F} \le q\delta_{1}\left\|X\right\|_{F},\tag{44}$$

for $\delta_1 = 4\sqrt{\mu^2 r^3 \log(2d)/(dq)}$.

Next, by the incoherence condition on X in equation (38), we have that

$$\|X\|_{\infty} \le \frac{4\mu(r+1)q}{\lambda}$$

Using equation (15) from Lemma 4.2 with $\nu = 4\mu(r+1)dq/\lambda$, we have

$$\|P_{\Omega}(XX^{\top})X - qXX^{\top}X\|_{F} \le \|P_{\Omega}(XX^{\top}) - qXX^{\top}\|_{2} \|X\|_{F} \le q\delta_{2} \|X\|_{F},$$
(45)

where $\delta_2 = 4\mu(r+1)\sqrt{dq\log(2d)}/\lambda \le 32\mu(r+1)\sqrt{\log(2d)/(dq)}$, for $\lambda \ge dq/8$.

We also need to add the error induced by N, given by

$$\|P_{\Omega}(N)X\|_{F} \leq \|P_{\Omega}(N)\|_{2} \cdot \|X\|_{F} \\ \leq q\delta_{3} \|X\|_{F}, \qquad (46)$$

where $\delta_3 = \sqrt{32\mu^2 r^2 \log(2d)/(dq)} + \tilde{O}(n^{-1}d^{-2})$, and the bound on the spectral norm of $P_{\Omega}(N)$ is shown in Corollary B.2.

Now, by plugging equations (44), (45), and (46) into the first-order optimality condition (cf. Eq. (36)), we have shown that

$$\left\| UU^{\top}X - XX^{\top}X - \gamma \nabla R(X) \right\|_{F} \le \delta \left\| X \right\|_{F}, \tag{47}$$

where $\delta = \delta_1 + \delta_2 + \delta_3$. When $||X||_F \leq \sigma_{\max}(U)\sqrt{r}$, both equations (42) and (43) are proved. Otherwise, by the triangle inequality

$$q \|UU^{\top}X\|_{F} \geq \|P_{\Omega}(UU^{\top})X\|_{F} - q\delta_{1} \|X\|_{F}$$

= $\|P_{\Omega}(XX^{\top})X + \lambda R(X)\|_{F} - q\delta_{1} \|X\|_{F} - \|P_{\Omega}(N)\|_{2} \|X\|_{F}$
$$\geq \|P_{\Omega}(XX^{\top})X\|_{F} - q(\delta_{1} + \delta_{3}) \|X\|_{F}$$

$$\geq q \|XX^{\top}X\|_{F} - q(\delta_{1} + \delta_{2} + \delta_{3}) \|X\|_{F}, \qquad (48)$$

where the second step is by equation (41) and the spectral norm bound on N, and the last step is again by the triangle inequality. Then, we have that

$$\left\| UU^{\top}X \right\|_{F}^{2} \le \left\| UU^{\top} \right\|_{2}^{2} \cdot \left\| X \right\|_{F}^{2} \le (\sigma_{\max}(U))^{4} \left\| X \right\|_{F}^{2}$$

On the other hand, $\|XX^{\top}X\|_{F}^{2} = \sum_{i=1}^{d} \sigma_{i}^{6}$, where σ_{i} is the *i*-th singular value of X, for i = 1, 2, ..., r. Notice that

$$\sum_{i=1}^{r} \sigma_i^2 \le \left(\sum_{i=1}^{r} \sigma_i^6\right)^{\frac{1}{3}} r^{\frac{2}{3}} \Rightarrow \left(\sum_{i=1}^{r} \sigma_i^2\right)^{\frac{1}{2}} \cdot r^{-1} \le \left(\sum_{i=1}^{r} \sigma_i^6\right)^{\frac{1}{2}}$$

by Hölder's inequality. Thus, plugging in this fact into the above equation (48), we obtain that

$$||X||_F^2 = \sum_{i=1}^r \sigma_i^2 \le r \cdot \sigma_{\max}(U))^2 + (\delta_1 + \delta_2 + \delta_3)r \le r \cdot (\sigma_{\max}(U))^2 + \delta_3$$

which implies that $||X||_F \leq r^{1/2}\sigma_{\max}(U) + \delta^{1/2}$. This concludes the proof of equation (42), which implies equation (43) based on equation (47).

Second, we look at the second-order optimality condition and show that this condition implies that the smallest singular value of X is lower bounded from below.

Lemma B.8. In the setting of Theorem 3.4, suppose X satisfies the second-order optimality condition (37) and equation (38). Suppose $\alpha \ge 4\kappa r \sqrt{\mu/d}$ and $\lambda = \mu(r+1)q/\alpha^2$. When $q \ge \frac{c\kappa^6 \mu^2 r^3 \log(d)}{d}$ for a fixed constant c > 0, and d is large enough, with probability at least $1 - O(d^{-1})$ over the randomness of Ω , we have that

$$(\sigma_{\min}(X))^2 \ge \frac{1}{8} (\sigma_{\min}(U))^2.$$
 (49)

Proof. Let $A = \{i : ||X_i|| \le \alpha\}$ be the index set of row vectors of X whose ℓ_2 norm is at most α . Let U_A be the matrix that has the same *i*-th row as the *i*-th row of U for every $i \in A$, and 0 elsewhere. Let $v \in \mathbb{R}^r$ be a unit vector such that $||Xv|| = \sigma_{\min}(X)$.

We show that

$$\sigma_{\min}(U_A) \ge (1 - (32\kappa)^{-1})\sigma_{\min}(U)$$

Let $B = \{1, 2, \dots, d\} \setminus A$. Since for any $i \in B$, $||X_i|| \ge \alpha$, we have that $|B|\alpha^2 \le ||X||_F^2 \le (\sigma_{\max}(U))^2 r + \delta$ by equation (42). It also follows that $|B| \le (\sigma_{\max}(U)^2 r + \delta)/\alpha^2$.

Next, we have that

$$\sigma_{\min}(U_A) = \sigma_{\min}(U - U_B) \ge \sigma_{\min}(U) - \sigma_{\max}(U_B)$$

$$\ge \sigma_{\min}(U) - \|U_B\|_F$$

$$\ge \sigma_{\min}(U) - \sqrt{|B|r\mu/d}$$

$$\ge \sigma_{\min}(U) - \sqrt{\frac{(\sigma_{\max}(U)^2 r + \delta)r\mu}{d\alpha^2}}$$

$$\ge (1 - \frac{1}{4})\sigma_{\min}(U), \qquad (50)$$

where the last line is because $\alpha \geq 4\kappa r \sqrt{\mu/d}$. It also follows that U_A has a column rank equal to r.

Notice that there must exist a unit vector u_A in the column span of U_A such that $||X^{\top}u_A|| \leq \sigma_{\min}(X)$. To see this, we can simply decompose the column span of U_A into one in the column span of X and another into the orthogonal subspace. Since u_A is a unit vector, we can also write u_A as $u_A = U_A\beta$, for some $\beta \in \mathbb{R}^r$ where $||\beta|| \leq \frac{1}{\sigma_{\min}(U_A)} \leq \frac{4}{3\sigma_{\min}(U)}$. Further,

$$\|u_A\|_{\infty} \le \left(\max_{i=1}^d \|U_A e_i\|\right) \cdot \|\beta\| \le \frac{4\sqrt{\mu r/d}}{3\sigma_{\min}(U)}$$

Next, we plug in $V = u_A v^{\top}$ in the second-order optimality condition of equation (37). Note that since u_A is in the column span of U_A , it is supported on the subset of columns spanned by A, and therefore $[\nabla^2 R(X)]V = 0.^8$ Therefore, the term about the Hessian of the regularization term in equation (37) is equal to zero. Thus, by taking $V = u_A v^{\top}$ in equation (37), we get that

$$\langle P_{\Omega}(u_A(Xv)^{\top} + Xvu_A^{\top}), u_A(Xv)^{\top} + Xvu_A^{\top} \rangle \ge 2u_A^{\top}P_{\Omega}(UU^{\top} - XX^{\top})u_A + 2u_A^{\top}Nu_A.$$
(51)

Note that

$$||Xv||_{\infty} \le (\max_{i=1}^{d} ||X_i||) ||v|| \le 2\sqrt{\mu(r+1)q/\lambda} = 2\alpha,$$

based on Lemma B.6 and the fact that v is a unit vector. As a result,

$$\left\| X v u_A^\top \right\|_{\infty} \le \frac{8\alpha \sqrt{\mu r/d}}{3\sigma_{\min}(U)} := \frac{\delta_1}{d},$$

⁸For details about the calculation regarding $\nabla^2 R(\cdot)$, see Lemma 18, Ge et al. (2017).

where $\delta_1 = \frac{8\kappa^2 r^{3/2} \mu}{3\sigma_{\min}(U)}$. By applying Proposition B.3 to the left-hand side of equation (51), we have that

$$\left| \langle P_{\Omega}(u_A(Xv)^{\top} + Xvu_A^{\top}), u_A(Xv)^{\top} + Xvu_A^{\top} \rangle - q \left\| u_A(Xv)^{\top} + Xvu_A^{\top} \right\|_F \right| \le q\delta_1 \sqrt{\frac{16\log(2d)}{dq}}.$$

Next, by applying Lemma 4.2 to the first term in the right-hand side of equation (51), we get

$$\left|u_A^\top P_\Omega(UU^\top - XX^\top)u_A - qu_A^\top(UU^\top - XX^\top)u_A\right| \le q\delta_2 \sqrt{\frac{16\log(2d)}{dq}},$$

where $\delta_2 = d \|UU^{\top} - XX^{\top}\|_{\infty} \leq \mu r + 4d\alpha^2 \leq 65\kappa\mu r$. Finally, we apply the spectral norm bound on N from equation (B.2) to $u_A^{\top}Nu_A$. Put together, we get

$$q \left\| X v u_A^\top + u_A (X v)^\top \right\|_F^2 \ge 2q \langle U U^\top - X X^\top, u_A u_A^\top \rangle - \delta q,$$

where

$$\delta = (\delta_1 + \delta_2) \sqrt{\frac{16\log(2d)}{dq}} + \sqrt{\frac{32\mu^2 r^2 \log(2d)}{dq}} + \tilde{O}(n^{-1/2}d^{-1}).$$

By simplifying the above calculation and using the fact that u_A is a unit vector, we get that

$$4 \|Xv\|^{2} + 2 \|X^{\top}u_{A}\|^{2} \ge 2 \|U^{\top}u_{A}\|^{2} - \delta.$$
(52)

Recall that u_A is a unit vector and u_A is in the column span of U_A . Therefore,

$$||U^{\top}u_A|| = ||U_A^{\top}u_A|| \ge (\sigma_{\min}(U_A))^2,$$

and the right hand side of equation (52) is greater than $2(\sigma_{\min}(U_A))^2 - \delta$.

Moreover, recall that $||Xv|| = \sigma_{\min}(X)$ and $||X^{\top}u_A|| \leq \sigma_{\min}(X)$, therefore, the left hand side of equation (52) is at most $6(\sigma_{\min}(X))^2$, in which we apply Cauchy-Schwarz to the cross term. Put together, from equation (52) we conclude that

$$(\sigma_{\min}(X))^2 \ge \frac{1}{3}(\sigma_{\min}(U_A))^2 - \frac{\delta}{6} \ge \frac{3}{16}(\sigma_{\min}(U))^2 - \frac{\delta}{6}$$

by equation (50). Thus, when $dq \ge c\kappa^6 r^3 \mu^3 \log(d)$, for $c = 16 \times 65^2 \times 9$, we have that $\frac{\delta}{6} \le (\sigma_{\min}(U))^2/16$, which concludes the proof of equation (49).

B.5 Characterization of The Population Loss

In the next result, we show that we could remove the effect of the regularizer term in the residual error. We use the fact that the regularizer follows the same direction as X, as stated in Proposition B.5.

Lemma B.9. In the setting of Theorem 3.4, suppose $X \in \mathbb{R}^{d \times r}$ satisfies that $(\sigma_{\min}(X))^2 \ge (\sigma_{\min}(U))^2/8$, and $\alpha \ge 4\kappa^2 r \sqrt{\mu/d}$,

$$\left\| UU^{\top}X - XX^{\top}X \right\|_{F}^{2} \leq \left\| UU^{\top}X - XX^{\top}X - \gamma\nabla R(X) \right\|_{F}^{2}, \text{ for any } \gamma \geq 0.$$

$$(53)$$

Proof. Let $S = \{i \in [d] : |X_i| \ge \alpha\}$ be the index set of row vectors of X whose ℓ_2 norm is greater than α . By the definition of the regularizer R(X), for any i not in S, we have that $(\nabla R(X))_i = 0$. Thus, for any $i \notin S$, we have that

$$\left\| U_i U^\top X - X_i X^\top X \right\|$$

For any i that is in S, we have that

$$\|U_{i}U^{\top}X - X_{i}X^{\top}X\|^{2} = \|U_{i}U^{\top}X - X_{i}X^{\top}X - (\gamma \nabla R(X))_{i}\|^{2},$$

where U_i, X_i are the *i*-th row vectors of $U, X \in \mathbb{R}^{d \times r}$, for every i = 1, 2, ..., d. As for each $i \in S$, we will show that

$$\langle (\nabla R(X))_i, X_i X^\top X - U_i U^\top X \rangle \ge 0.$$
(54)

By Proposition B.5, we have $(\nabla R(X))_i = \left(4\lambda \sum_{i=1}^d \frac{(|X_i|^3 - \alpha)_+^3}{|X_i|}\right) X_i := a_i X_i$, where $a_i \ge 0$ denotes the multiplier in front of X_i . Then, we have that

$$\langle (\nabla R(X))_i, X_i X^\top X \rangle = a_i \left\| X X_i^\top \right\|^2 \ge a_i \left\| X_i \right\|^2 \sigma_{\min} \left(X^\top X \right)$$
$$\ge \frac{a_i}{8} \left\| X_i \right\|^2 (\sigma_{\min}(U))^2,$$

by the condition on the smallest singular value of U. On the other hand,

$$\begin{aligned} \langle (\nabla R(X))_i, U_i U^\top X \rangle &= a_i \langle X_i, U_i U^\top X \rangle \\ &\leq a_i \|X_i\| \cdot \|U_i\| \cdot \sigma_{\max}(U) \cdot \sigma_{\max}(X) \\ &\leq a_i \|X_i\| \|U_i\| (\sigma_{\max}(U))^2 2\sqrt{r} \qquad \text{(by equation (42), plus } \sigma_{\max}(X) \leq \|X\|_F) \\ &\leq a_i \|X_i\| \|U_i\| (\sigma_{\min}(U))^2 2\kappa^2 \sqrt{r} \qquad (\text{since } \sigma_{\max}(U) = \kappa \cdot \sigma_{\min}(U)) \\ &\leq \frac{a_i}{16} \|X_i\|^2 (\sigma_{\min}(U))^2, \end{aligned}$$

where the last step is because $||X_i|| \ge \alpha \ge 4\kappa^2 r \sqrt{\frac{\mu}{d}} \ge 4\kappa^2 \sqrt{r} ||U_i||$, since $||U_i|| \le \sqrt{\mu r/d}$ by Assumption 3.3. Putting both sides together, we therefore conclude that equation (54) must hold. The proof is completed.

Finally, we show that the above result also implies that UU^{\top} is close to XX^{\top} .

Lemma B.10. In the setting of Theorem 3.4, suppose $X \in \mathbb{R}^{d \times r}$ satisfies equations (43) and (53). Then, with probability at least $1 - O(d^{-1})$, the following must be true

$$\|XX^{\top} - UU^{\top}\|_{F}^{2} \le \frac{c\mu^{2}r^{5}\kappa^{6}\log(2d)}{dq},$$
(55)

for a fixed constant c > 0.

Proof. We separate U into one part that is in the column span of X and another part that is orthogonal to the column span of X. Let U = Z + V where Z is the projection of U to the column span of X, and V is the projection of U to the orthogonal column subspace of X. In particular, we have $V^{\top}X = 0$. Thus,

$$UU^{\top}X = (Z+V)(Z+V)^{\top}X = ZZ^{\top}X + VZ^{\top}X$$

$$\Rightarrow \left\|UU^{\top}X - XX^{\top}X\right\|^{2} = \left\|ZZ^{\top}X - XX^{\top}X\right\|^{2}_{F} + \left\|VZ^{\top}X\right\|^{2}_{F} \le \delta^{2},$$
(56)

for $\delta = 42\sigma_{\max}(U)\sqrt{\frac{\mu^2 r^4 \log(2d)}{dq}}$, by equation (43) and then equation (53).

Notice that Z has the same column span as X. As a result,

$$\left\| ZZ^{\top}X - XX^{\top}X \right\|_{F}^{2} \ge \sigma_{\min}^{2}(X) \left\| ZZ^{\top} - XX^{\top} \right\|_{F}^{2},$$

which can be verified by inserting the SVD of X into the above inequality. Combined with equation (56), we thus have that

$$\left\| ZZ^{\top} - XX^{\top} \right\|_F^2 \le \frac{\delta^2}{(\sigma_{\min}(X))^2} \le \frac{8\delta^2}{(\sigma_{\min}(U))^2},$$

since $(\sigma_{\min}(X))^2 \ge (\sigma_{\min}(U))^2/8$.

Next, let $a \in \mathbb{R}^d, b \in \mathbb{R}^r$ be two unit vectors such that $a^\top Z Z^\top X b = \sigma_{\min}(Z Z^\top X)$. From equation (56) we can infer that

$$\left|\sigma_{\min}(ZZ^{\top}X) - a^{\top}XX^{\top}Xb\right| \le \delta.$$
(57)

Now we consider two cases. Suppose $a^{\top}XX^{\top}Xb$ is positive. Then the left-hand side of equation (57) must be at least

$$a^{\top}XX^{\top}Xb - \sigma_{\min}(ZZ^{\top}X) = \sigma_{\min}(XX^{\top}X) - \sigma_{\min}(ZZ^{\top}X).$$

The other case is that if $a^{\top}XX^{\top}Xb$ is negative, then the left hand side of equation (57) must be at least

$$\left|a^{\top}XX^{\top}Xb\right| - \sigma_{\min}(ZZ^{\top}X) \ge \sigma_{\min}(XX^{\top}X) - \sigma_{\min}(ZZ^{\top}X)$$

We know that $\sigma_{\min}(XX^{\top}X) \geq (\sigma_{\min}(U))^3/(8\sqrt{8})$. Thus, if $\delta \leq (\sigma_{\min}(U))^3/(16\sqrt{8})$, then

$$\sigma_{\min}(ZZ^{+}X) \ge (\sigma_{\min}(U))^{3}/(16\sqrt{8}).$$

On the other hand,

$$\sigma_{\min}(ZZ^{\top}X) \le \sigma_{\max}(Z) \cdot \sigma_{\min}(Z^{\top}X) \le \|U\|_2 \cdot \sigma_{\min}(Z^{\top}X),$$

which implies

$$\sigma_{\min}(Z^{\top}X) \ge \frac{(\sigma_{\min}(U))^3}{8\sqrt{8} \|U\|_2}.$$

Also note that $\left\| V Z^{\top} X \right\|_{F}^{2} \leq \delta^{2}$, which implies that

$$\|V\|_F^2 \le \frac{\delta^2}{(\sigma_{\min}(Z^\top X))^2} \le \frac{512\delta^2 \|U\|_2^2}{(\sigma_{\min}(U))^6} = \frac{512\delta^2 \kappa^2}{(\sigma_{\min}(U))^4}$$

Finally, $||ZV^{\top}||_F \leq ||Z||_2 \cdot ||V||_F \leq ||U||_2 \cdot ||V||_F$, since Z is a projection of U to the column space of X. In summary, we have that

$$\begin{split} \left\| UU^{\top} - XX^{\top} \right\|_{F}^{2} &= \left\| ZZ^{\top} - XX^{\top} \right\|_{F}^{2} + 2 \left\| ZV^{\top} \right\|_{F}^{2} + \left\| VV^{\top} \right\|_{F}^{2} \\ &\leq \frac{16\delta^{2}}{(\sigma_{\min}(U))^{2}} + \frac{1024\delta^{2}\kappa^{4}}{(\sigma_{\min}(U))^{2}} + \frac{512^{2}\delta^{4}\kappa^{4}r}{(\sigma_{\min}(U))^{8}} \\ &\leq \frac{(512\kappa^{4}r + 1024\kappa^{4} + 16)\delta^{2}}{(\sigma_{\min}(U))^{2}}, \end{split}$$

since $\delta \leq (\sigma_{\min}(U))^3/(8\sqrt{8})$, which concludes the proof of equation (55), for $c = 42^2 \times 512$.

Now we are ready to complete the proof of Theorem 3.4.

Proof of Theorem 3.4. Suppose X satisfies the first and second order optimality conditions. Then by Lemma B.6 and Lemma B.7, we have that X satisfies the incoherence condition in equation (38). As a result, $(\sigma_{\min}(X))^2 \ge (\sigma_{\min}(U))^2/8$ by Lemma B.8. In particular, we require α to be at least $4\kappa^2 r \sqrt{\mu/d}$. In addition, based on Lemma B.7, we require $\alpha \le \sqrt{\mu(r+1)q/\lambda}$, which requires $\lambda \ge \mu(r+1)q/\alpha^2 \ge dq/8$.

Recall that $q = 1 - (1 - p^2)^n = 1 - (1 - C^2/d^2)^n \approx C^2 n/d^2$. When $n \ge \frac{cdr^5 \kappa^6 \mu^2 \log(2d)}{C^2 \epsilon^2}$, the right hand side of equation (55) is at most ϵ^2 . This concludes the proof of equation (10).

Remark B.11. Our proof has been inspired by the important early work of Ge et al. (2016). We complement their work in three aspects.

- First, in the one-sided completion setting, the sampling patterns are non-random. In particular, the diagonal entries are fully observed, while the off-diagonal entries are sparsely observed. Thus, there is a mix of dense and sparse samplings in Ω . To address this problem, we use a weighted projection as defined in Section 2.
- Second, the bias of Hájek's estimator is complex due to the nonlinear form of the estimator, and requires a delicate analysis of the noise matrix N. This is studied in detail in Corollary B.1 and Corollary B.2.
- Third, the observation patterns in the off-diagonal entries of Ω are not completely independent, requiring a new, row-by-row concentration inequality that leverages a conditional independence property, i.e., Lemma 4.2. More broadly, non-uniform sampling patterns are very common in real-world datasets and have received lots of recent interest in causal inference (Athey et al., 2021; Xiong & Pelger, 2023). We hope this work inspires further studies of modeling non-uniform sampling on sparse panel data from an optimization perspective.

C Omitted Experiments

First, we provide the full procedure for recovering the full matrix using one-sided matrix completion in Algorithm 2, corresponding to Figure 3.

Algorithm 2 Missing data imputation using one-sided matrix completion

Input: A partially observed $\hat{M} \in \mathbb{R}^{n \times d}$

Require: Rank r, number of iterations t, learning rate η

Output: An k by d matrix corresponding to the imputed rows of S

1: $Z \leftarrow H\acute{A}$ JEK-GD $(\hat{M}; r, t, \eta)$ // May add Gaussian noise to the non-zero entries of \hat{M} 2: $U_r D_r U_r^\top \leftarrow \text{Rank-}r$ SVD of Z

3: $\Omega \leftarrow$ The set of indices corresponding to the nonzero entries of $\hat{M}^{\top}\hat{M}$

4:
$$Q \leftarrow \arg\min_{Q \in \mathbb{R}^{n \times r}} \frac{1}{2} \sum_{i \in S, (i,j) \in \Omega} \left(\left(QU_r^\top \right)_{i,j} - \hat{M}_{i,j} \right)$$

5: return QX^{\top}

Next, we describe dataset statistics and pre-processing procedures. The MovieLens-32M dataset comprises 32 million user ratings on approximately 87,000 movies, collected from the MovieLens recommendation platform, with ratings ranging from 1 to 5. The sparsity ratio is roughly 3×10^{-4} . We also consider two related MovieLens datasets with 20 million and 25 million entries to ensure the robustness of our findings. We split 80% into training and hold out the rest for testing.

The Amazon Reviews dataset includes 571 million user ratings on 48 million products, collected from Amazon, with ratings ranging from 1 to 5. The sparsity ratio is roughly 2×10^{-7} . For the Amazon reviews dataset, we choose the Automotive category with 50,000 items rated by 100,000 users as M. We use 80-20 split.

The Genomes dataset contains phased genotype data for 2,054 individuals gathered from diverse populations. The dataset is represented as a dense matrix, with rows representing genomic sites and columns corresponding to individuals. Each entry represents a biallelic genotype, encoded as 0 or 1. We map them to 1 and 2, and use 0 to indicate missing entries. This dataset is fully observed. Thus, we use 1% for training and hold out the rest for testing.

The statistics of the five real-world datasets, represented as sparse matrices, are summarized in Table 3 below. For MovieLens and Amazon Reviews datasets, we use \hat{T} on $M^{\top}M$ to obtain the ground truth second-moment matrix. Then, the estimation error is evaluated against this matrix. In our experiments, we set the sampling probability p = 0.8. Given that the original datasets have an average sparsity of approximately 0.3%, even with p = 0.8, the number of observed entries remains very sparse. When T is not fully observed, we focus on measuring the mean squared error within the observed entries instead and treat the missing entries as zero.

Dataset	Genomes	Amazon Reviews	MovieLens-20M	MovieLens-25M	MovieLens-32M
# rows # columns # nonzero entries	$ \begin{array}{c c} 100,000 \\ 2,504 \\ 2.5 \times 10^8 \end{array} $	$egin{array}{c} 100,000\ 50,000\ 1.3 imes 10^6 \end{array}$	138,000 27,000 2×10^{7}	162,000 62,000 2.5×10^{7}	200,948 87,585 3.2×10^7

Table 3: We describe the detailed statistics of five real-world datasets used in our experiments.

Implementation of baseline methods. We provide a detailed description of all the baseline methods we have implemented in our experiments. Alternating-GD minimizes the squared reconstruction error on the observed entries of M. We minimize the following objective:

$$\min_{X \in \mathbb{R}^{n \times r}, Y \in \mathbb{R}^{d \times r}} \|P_{\Omega}(XY^{\top} - M)\|_{F}^{2}$$

where P_{Ω} denotes the projection onto the observed entries. We initialize X and Y from an isotropic Gaussian distribution. We perform the optimization using Adam with a learning rate of 0.1 over 300 training epochs.

SoftImpute-ALS alternates between matrix completion using singular value decomposition (SVD) and imputing missing values (Hastie et al., 2015). At each iteration, missing entries are filled with the current estimates, followed by SVD and soft thresholding of the singular values. The resulting truncated low-rank approximation is used to update the matrix. We repeat this process until convergence, as measured by the Frobenius norm difference between successive reconstructions, or until a maximum of 500 iterations. We set the target rank to the same as r and the convergence threshold to 10^{-5} .



Figure 6: Recovery error of T (left two columns) and M (right two columns) using our approach on synthetic data. Figures 6a, 6b: One-sided recovery error under uniform sampling by varying p, and sampling two entries per row (C = 2) while varying the number of rows n. Figures 6c, 6d: Recovery error on M, also under uniform sampling, and sampling two entries per row. Overall, we find that our approach consistently provides more accurate estimates compared with the three baseline methods, all of which are widely used in the matrix completion literature. Figures 6e, 6f, 6g, 6h: Vary C and repeat the same experiment.

Table 4: We report the recovery error for both one-sided matrix completion (i.e., recovering T), and recoverin	g
M, along with the corresponding running time (measured in seconds on an Ubuntu server) on three MovieLen	ıs
datasets. We run each experiment with five random seeds and report the mean and standard deviation.	

Dataset	MovieLens-20M	MovieLens-25M	MovieLens-20M	MovieLens-25M
T	Recovery error		Running time (Seconds)	
Alternating-GD SoftImpute-ALS Algorithm 1	$\begin{array}{c} 0.015_{\pm 0.009} \\ 0.008_{\pm 0.005} \\ \textbf{0.002}_{\pm 0.001} \end{array}$	$\begin{array}{c} 0.006_{\pm 0.002} \\ 0.006_{\pm 0.002} \\ \textbf{0.002}_{\pm 0.10} \end{array}$	$ \begin{vmatrix} 1.9_{\pm 0.1} \times 10^2 \\ 3.4_{\pm 0.1} \times 10^3 \\ 4.8_{\pm 0.2} \times \mathbf{10^1} \end{vmatrix} $	$\begin{array}{c} 4.4_{\pm 0.1} \times 10^2 \\ 9.4_{\pm 0.7} \times 10^3 \\ 1.1_{\pm 0.6} \times \mathbf{10^2} \end{array}$
M	Root mean squared error		Running tim	e (Seconds)
Alternating-GD SoftImpute-ALS Algorithm 2	$\begin{array}{c c} 1.11_{\pm 0.00} \\ 1.09_{\pm 0.00} \\ 0.99_{\pm 0.00} \end{array}$	$\begin{array}{c} 1.27_{\pm 0.02} \\ 1.25_{\pm 0.00} \\ \textbf{1.04}_{\pm \textbf{0.00}} \end{array}$	$ \begin{vmatrix} 1.9_{\pm 0.1} \times 10^2 \\ 3.4_{\pm 0.1} \times 10^3 \\ 1.0_{\pm 0.1} \times 10^2 \end{vmatrix} $	$\begin{array}{c} 4.4_{\pm 0.1} \times 10^2 \\ 9.4_{\pm 0.7} \times 10^3 \\ \textbf{2.6}_{\pm 0.1} \times \textbf{10^2} \end{array}$

The nuclear norm regularization approach is another commonly used approach in matrix completion (Zhang et al., 2019a). In one-sided estimation, we compute empirical averages based on the observed data (Cao et al., 2023). For off-diagonal entries (i, j), we identify all samples in which the pair (i, j) appears and average the product of the corresponding entries from the two relevant columns. For diagonal entries (i, i), we consider all samples where index i appears in either position of a pair, square the corresponding data matrix entry for each such sample, and average these squared values. We then estimate the full matrix using symmetric alternating gradient descent by solving the following optimization problem:

$$\min_{X \in \mathbb{R}^{d \times r}} \left\| P_{\Omega}(XX^{\top} - T) \right\|_{F}^{2} + \lambda \left\| X \right\|_{\star},$$

where $||X||_{\star}$ denotes the nuclear norm of X and λ is set as 0.01. We perform the optimization using Adam with a learning rate of 0.1 for 1,000 steps.

Hyperparameter configurations. We report the hyperparameter configurations used in our experiments. The learning rate η is varied from 10 to 10^5 , the regularization coefficient λ from 10^{-2} to 10^{-5} , and the regularization parameter α from 10^{-4} to 1. Based on cross-validation, we select $\eta = 10^4$, $\lambda = 10^{-4}$, and $\alpha = 10^{-3}$, which yield the lowest error. For the Amazon Review dataset, we set $\lambda = 0.01$ and $\alpha = 0.1$.

C.1 Comparison Results

First, we report the runtime comparison between HÁJEK-GD and nuclear norm regularization. We run gradient descent with 1,000 iterations until it has converged. We also fix the number of observed entries as 200*d* while varying *d* from 10^3 to 10^4 . We find that our approach requires 0.95 seconds for $d = 10^3$ and 1.61 seconds for $d = 10^4$. By contrast, solving a convex program with nuclear norm penalty takes 5.68 seconds for $d = 10^3$ and 100.34 seconds for $d = 10^4$, measured on the server.

Second, we present the recovery error for imputing M, illustrated in Figures 6. Our approach consistently reduces the recovery error compared to all three baseline methods. For instance, when C = 2 (sampling two entries per row), our approach incurs an error of 0.0003, compared to 0.002 for nuclear norm regularization. We also observe similar results on synthetic data when sampling five entries per row (C = 5) or sampling ten entries per row (C = 10) in Figures 6e and 6f. Similar results are observed when our approach is applied to recover the missing entries of M. For recovering M, our approach reduces estimation error by 94% relative to alternating-GD and by 98% compared to softImpute-ALS.

Finally, we report the results on MovieLens-20M and MovieLens-25M in Table 4. For one-sided matrix completion, our approach reduces the Frobenius norm error between the estimate and the true T by up to 87% compared to alternating-GD and up to 75% compared to softImpute-ALS, while also achieving lower running times. For recovering the missing entries of M, our approach reduces the RMSE by up to 10% compared to alternating-GD and up to 6% compared to softImpute-ALS.

C.2 Analysis on Rank and Sensitivity

We use a target rank of 10 for all experiments reported in Section 5. To further assess the robustness of the method, we conduct additional experiments on three real-world datasets in Algorithm 1, with the target rank in the range of 1 to 30. The results are shown in Table 5. We find that larger ranks may lead to overfitting in sparse settings; for example, on the Amazon dataset, increasing r beyond 10 results in worse performance when recovering M from the estimated T.

We also run similar analyses on synthetic data and find that the algorithm is not particularly sensitive to different choices in this context. In particular, we conduct different simulations by varying the rank r of M between 1 and 50. We vary the target rank used in Algorithm 1 between 1 and 50. We find that the estimation errors are comparable between different choices of target ranks.

Table 5: We vary the rank of the variable matrix in HÁJEK-GD, corresponding to the results we reported in Table 2. The results are averaged over five independent runs.

Dataset	MovieLens-20M	Amazon Reviews	Genomes		
T (Using Algorithm 1)	One-sided recovery error				
r = 1	$7.0_{\pm 0.1} \times 10^{-3}$	$4.5_{\pm 0.1} \times 10^{-3}$	$5.8_{\pm 0.1} \times 10^{-5}$		
r = 10	$4.7_{\pm 0.1} \times 10^{-3}$	$1.3_{\pm 0.1} imes 10^{-3}$	$5.8_{\pm0.1} imes 10^{-5}$		
r = 20	$1.9_{\pm 0.1} \times 10^{-3}$	$1.3_{\pm 0.1} \times 10^{-3}$	$5.8_{\pm0.1} \times 10^{-5}$		
r = 30	$1.8_{\pm 0.1} imes 10^{-3}$	$1.3_{\pm 0.1} \times 10^{-3}$	$5.8_{\pm 0.1} \times 10^{-5}$		
M (Using Algorithm 2)	Root mean squared recovery error				
r = 1	$3.7_{\pm 0.1}$	$4.6_{\pm 0.1}$	$0.1_{\pm 0.1}$		
r = 10	$1.1_{\pm 0.1}$	$1.9_{\pm 0.1}$	$0.2_{\pm 0.1}$		
r = 20	$1.0_{\pm 0.1}$	$2.3_{\pm 0.1}$	$0.2_{\pm 0.1}$		
r = 30	$1.0_{\pm 0.1}$	$2.7_{\pm 0.1}$	$0.2_{\pm 0.1}$		

Sensitivity of Algorithm 1. We evaluate the sensitivity of adding noise to \hat{M} on the final output of Algorithm 1. We use synthetic datasets with varied numbers of rows between 5d and 20d, and a fixed dimension of 1,000. We also test different ranks of M, ranging from 5 to 20 We add Gaussian noise with a standard deviation that ranges from 0 to 0.01 to the input. Interestingly, we find that estimation error increases linearly with σ . Moreover, when varying n, the sensitivity level of HÁJEK-GD (recall its definition in equation (17)) is generally low. When n = 5d, the slope of the line is 19. When n = 10d, the slope decreases to 14. When n = 20d, the slope further decreases to 13.

Similar results hold after varying the rank between 20, 10, 5, and the sensitivity level ranges between 16, 13, 11. See illustration in Figure 7.



(a) Varying noise level σ for different number of rows n

(b) Varying noise level σ for different rank r of M

Figure 7: Illustrating the estimation error after adding Gaussian noise with mean zero and variance σ^2 to \hat{M} on the non-zero entries. In Figure 7a, we vary the number of rows *n* with a fixed rank r = 10. In Figure 7b, we vary rank *r* (of *M*) with a fixed m = 20d. Across all six cases, the sensitivity level, measured by the slope of the line, is at most 19 and drops to 11 as the sample size increases or the rank decreases. We report the mean and standard deviation from five independent runs.