

Regularization for Matrix Completion

Raghunandan H. Keshavan* and Andrea Montanari*[†]
 Departments of Electrical Engineering* and Statistics[†], Stanford University

Abstract—We consider the problem of reconstructing a low rank matrix from noisy observations of a subset of its entries. This task has applications in statistical learning, computer vision, and signal processing. In these contexts, ‘noise’ generically refers to any contribution to the data that is not captured by the low-rank model. In most applications, the noise level is large compared to the underlying signal and it is important to avoid overfitting. In order to tackle this problem, we define a *regularized* cost function well suited for spectral reconstruction methods. Within a random noise model, and in the large system limit, we prove that the resulting accuracy undergoes a phase transition depending on the noise level and on the fraction of observed entries. The cost function can be minimized using OPTSPACE (a manifold gradient descent algorithm). Numerical simulations show that this approach is competitive with state-of-the-art alternatives.

I. INTRODUCTION

Let N be an $m \times n$ matrix which is ‘approximately’ low rank, that is

$$N = M + W = U\Sigma V^T + W. \quad (1)$$

where U has dimensions $m \times r$, V has dimensions $n \times r$, and Σ is a diagonal $r \times r$ matrix. Thus M has rank r and W can be thought of as noise, or ‘unexplained contributions’ to N . Throughout the paper we assume the normalization $U^T U = m I_{r \times r}$ and $V^T V = n I_{r \times r}$ ($I_{d \times d}$ being the $d \times d$ identity).

Out of the $m \times n$ entries of N , a subset $E \subseteq [m] \times [n]$ is observed. We let $\mathcal{P}_E(N)$ be the $m \times n$ matrix that contains the observed entries of N , and is filled with 0’s in the other positions

$$\mathcal{P}_E(N)_{ij} = \begin{cases} N_{ij} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

The *noisy matrix completion* problem requires to reconstruct the low rank matrix M from the observations $\mathcal{P}_E(N)$. In the following we will also write $N^E = \mathcal{P}_E(N)$ for the sparsified matrix. Over the last year, matrix completion has attracted significant attention because of its relevance –among other applications– to collaborative filtering. In this case, the matrix N contains evaluations of a group of customers on a group of products, and one is interested in exploiting a sparsely filled matrix to provide personalized recommendations [1].

In such applications, the noise W is not a small perturbation and it is crucial to avoid overfitting. For instance, in the limit $M \rightarrow 0$, the estimate of \widehat{M} risks to be a low-rank approximation of the noise W , which would be grossly incorrect.

In order to overcome this problem, we propose in this paper an algorithm based on minimizing the following cost function

$$\mathcal{F}_E(X, Y; S) \equiv \frac{1}{2} \|\mathcal{P}_E(N - XSY^T)\|_F^2 + \frac{1}{2} \lambda \|S\|_F^2. \quad (3)$$

Here the minimization variables are $S \in \mathbb{R}^{r \times r}$, and $X \in \mathbb{R}^{m \times r}$, $Y \in \mathbb{R}^{n \times r}$ with $X^T X = Y^T Y = I_{r \times r}$. Finally, $\lambda > 0$ is a regularization parameter.

A. Algorithm and main results

The algorithm is an adaptation of the OPTSPACE algorithm developed in [2]. A key observation is that the following modified cost function can be minimized by singular value decomposition (see Section I.1):

$$\widehat{\mathcal{F}}_E(X, Y; S) \equiv \frac{1}{2} \|\mathcal{P}_E(N) - XSY^T\|_F^2 + \frac{1}{2} \lambda \|S\|_F^2. \quad (4)$$

As emphasized in [2], [3], which analyzed the case $\lambda = 0$, this minimization can yield poor results unless the set of observations E is ‘well balanced’. This problem can be bypassed by ‘trimming’ the set E , and constructing a balanced set \widetilde{E} . The OPTSPACE algorithm is given as follows.

OPTSPACE (set E , matrix N^E)

- 1: Trim E , and let \widetilde{E} be the output;
 - 2: Minimize $\widehat{\mathcal{F}}_{\widetilde{E}}(X, Y; S)$ via SVD, let X_0, Y_0, S_0 be the output;
 - 3: Minimize $\mathcal{F}_E(X, Y; S)$ by gradient descent using X_0, Y_0, S_0 as initial condition.
-

In this paper we will study this algorithm under a model for which step 1 (trimming) is never called, i.e. $\widetilde{E} = E$ with high probability. We will therefore not discuss it any further. Section II compares the behavior of the present approach with alternative schemes. Our main analytical result is a sharp characterization of the mean square error after step 2. Here and below the limit $n \rightarrow \infty$ is understood to be taken with $m/n \rightarrow \alpha \in (0, \infty)$.

Theorem I.1. *Assume $|M_{ij}| \leq M_{\max}$, W_{ij} to be i.i.d. random variables with mean 0 variance $\sqrt{mn}\sigma^2$ and $\mathbb{E}\{W_{ij}^4\} \leq Cn^2$, and that for each entry (i, j) , N_{ij} is observed (i.e. $(i, j) \in E$) independently with probability p . Finally let $\widehat{M} = X_0 S_0 Y_0^T$ be the rank r matrix reconstructed by step 2 of OPTSPACE, for the optimal choice of λ . Then, almost surely for $n \rightarrow \infty$*

$$\frac{1}{\|\widehat{M}\|_F^2} \|\widehat{M} - M\|_F^2 = 1 - \frac{\left\{ \sum_{k=1}^r \Sigma_k^2 \left(1 - \frac{\sigma^4}{p^2 \Sigma_k^4} \right) \right\}_+^2}{\|\Sigma\|_F^2 \left\{ \sum_{k=1}^r \Sigma_k^2 \left(1 + \frac{\sqrt{\alpha} \sigma^2}{p \Sigma_k^2} \right) \left(1 + \frac{\sigma^2}{p \Sigma_k^2 \sqrt{\alpha}} \right) \right\}} + o_n(1).$$

This theorem focuses on a high-noise regime, and predicts a sharp phase transition: if $\sigma^2/p < \Sigma_1$, we can successfully extract information on M , from the observations N^E . If on

the other hand $\sigma^2/p \geq \Sigma_1$, the observations are essentially useless in reconstructing M . It is possible to prove [4] that the resulting tradeoff between noise and observed entries is tight: no algorithm can obtain relative mean square error smaller than one for $\sigma^2/p \geq \Sigma_1$, under a simple random model for M . To the best of our knowledge, this is the first sharp phase transition result for low rank matrix completion.

For the proof of Theorem I.1, we refer to Section III. An important byproduct of the proof is that it provides a rule for choosing the regularization parameter λ , in the large system limit.

B. Related work

The importance of regularization in matrix completion is well known to practitioners. For instance, one important component of many algorithms competing for the Netflix challenge [1], consisted in minimizing the cost function $\mathcal{H}_E(X, Y; S) \equiv \frac{1}{2} \|\mathcal{P}_E(N - \tilde{X}\tilde{Y}^T)\|_F^2 + \frac{1}{2} \lambda \|\tilde{X}\|_F^2 + \frac{1}{2} \lambda \|\tilde{Y}\|_F^2$ (this is also known as *maximum margin matrix factorization* [5], [6]). Here the minimization variables are $\tilde{X} \in \mathbb{R}^{m \times r}$, $\tilde{Y} \in \mathbb{R}^{n \times r}$. Unlike in OPTSPACE, these matrices are not constrained to be orthogonal, and as a consequence the problem becomes significantly more degenerate. Notice that, in our approach, the orthogonality constraint fixes the norms $\|\tilde{X}\|_F$, $\|\tilde{Y}\|_F$. This motivates the use of $\|S\|_F^2$ as a regularization term.

Convex relaxations of the matrix completion problem were recently studied in [7], [8]. As emphasized by Mazumder, Hastie and Tibshirani [9], such nuclear norms relaxations can be viewed as spectral regularizations of a least square problem. Finally, the phase transition phenomenon in Theorem I.1, generalizes a result of Johnstone and Lu on principal component analysis [10], and similar random matrix models were studied in [11].

II. NUMERICAL SIMULATIONS

In this section, we present the results of numerical simulations on synthetically generated matrices. The data are generated following the recipe of [9]: sample $\bar{U} \in \mathbb{R}^{n \times r}$ and $\bar{V} \in \mathbb{R}^{m \times r}$ by choosing \bar{U}_{ij} and \bar{V}_{ij} independently and identically as $\mathcal{N}(0, 1)$. Sample independently $W \in \mathbb{R}^{m \times n}$ by choosing W_{ij} iid with distribution $\mathcal{N}(0, \sigma^2 \sqrt{mn})$. Set $N = \bar{U}\bar{V}^T + W$. We also use the parameters chosen in [9] and define

$$\begin{aligned} \text{SNR} &= \sqrt{\frac{\text{Var}((\bar{U}\bar{V}^T)_{ij})}{\text{Var}(W_{ij})}}, \\ \text{TestError} &= \frac{\|\mathcal{P}_E^\perp(\bar{U}\bar{V}^T - \hat{N})\|_F^2}{\|\mathcal{P}_E^\perp(\bar{U}\bar{V}^T)\|_F^2}, \\ \text{TrainError} &= \frac{\|\mathcal{P}_E(N - \hat{N})\|_F^2}{\|\mathcal{P}_E(N)\|_F^2}, \end{aligned}$$

where $\mathcal{P}_E^\perp(A) \equiv A - \mathcal{P}_E(A)$.

In Figure 1, we plot the train error and test error for the OPTSPACE algorithm on matrices generated as above with $n = 100, r = 10, \text{SNR}=1$ and $p = 0.5$. For comparison, we

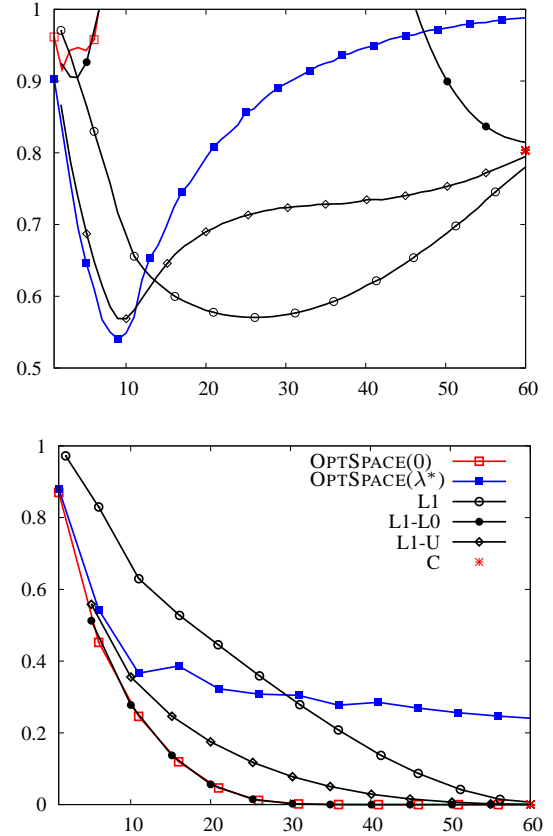


Fig. 1. Test (top) and train (bottom) error vs. rank for OPTSPACE, SOFT-IMPUTE, HARD-IMPUTE and SVT. Here $m = n = 100, r = 10, p = 0.5, \text{SNR} = 1$.

also plot the corresponding curves for SOFT-IMPUTE, HARD-IMPUTE and SVT taken from [9]. In Figures 2 and 3, we plot the same curves for different values of r, ϵ, SNR . In these plots, OPTSPACE(λ) corresponds to the algorithm that minimizes the cost (3). In particular OPTSPACE(0) corresponds to the algorithm described in [2]. Further, $\lambda^* = \lambda^*(\rho)$ is the value of the regularization parameter that minimizes the test error while using rank ρ (this can be estimated on a subset of the data, not used for training).

It is clear that regularization greatly improves the performance of OPTSPACE and makes it competitive with the best alternative methods.

III. PROOF OF THEOREM I.1

The proof of Theorem 1 is based on the following three steps: (i) Obtain an explicit expression for the root mean square error in terms of right and left singular vectors of N ; (ii) Estimate the effect of the noise W on the right and left singular vectors; (iii) Estimate the effect of missing entries. Step (ii) builds on recent estimates on the eigenvectors of large covariance matrices [12]. In step (iii) we use the results of [2]. Step (i) is based on the following linear algebra calculation, whose proof we omit due to space constraints (here and below $\langle A, B \rangle \equiv \text{Tr}(AB^T)$).

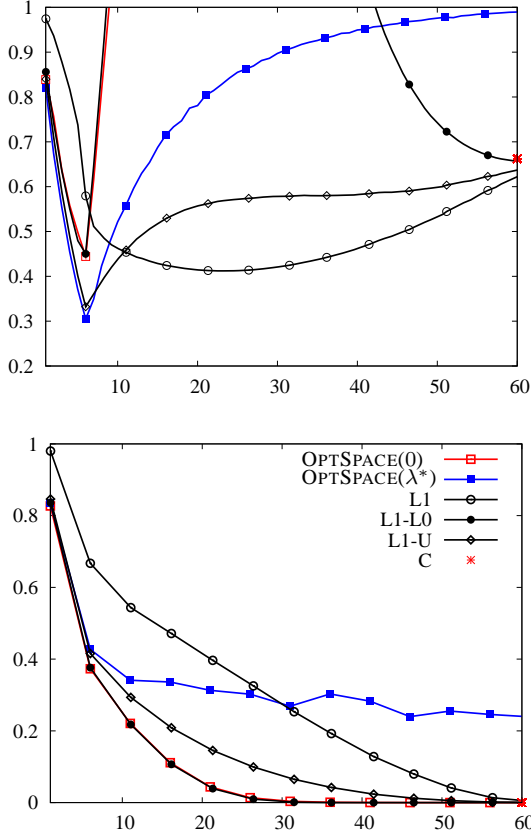


Fig. 2. Test (top) and train (bottom) error vs. rank for OPTSPACE, SOFT-IMPUTE, HARD-IMPUTE and SVT. Here $m = n = 100, r = 6, p = 0.5, \text{SNR} = 1$.

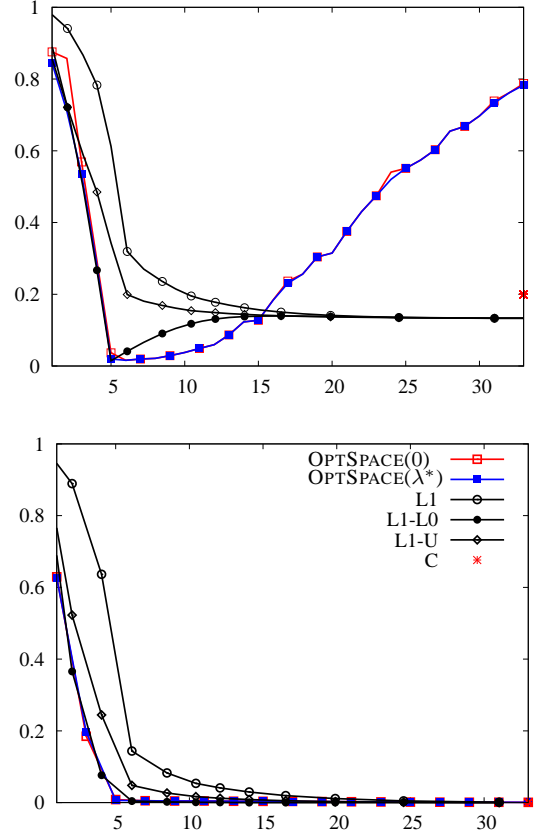


Fig. 3. Test (top) and train (bottom) error vs. rank for OPTSPACE, SOFT-IMPUTE, HARD-IMPUTE and SVT. Here $m = n = 100, r = 5, p = 0.2, \text{SNR} = 10$.

Proposition III.1. Let $X_0 \in \mathbb{R}^{m \times r}$ and $Y_0 \in \mathbb{R}^{m \times r}$ be the matrices whose columns are the first r , right and left, singular vectors of N^E . Then the rank- r matrix reconstructed by step 2 of of OPTSPACE, with regularization parameter λ , has the form $\widehat{M}(\lambda) = X_0 S_0(\lambda) Y_0^T$. Further, there exists $\lambda_* > 0$ such that

$$\frac{1}{mn} \|M - \widehat{M}(\lambda_*)\|_F^2 = \|\Sigma\|_F^2 - \left(\frac{\langle X_0^T M Y_0, X_0^T N^E Y_0 \rangle}{\sqrt{mn} \|X_0 N^E Y_0\|_F} \right)^2. \quad (5)$$

A. The effect of noise

In order to isolate the effect of noise, we consider the matrix $\widehat{N} = p U \Sigma V^T + W^E$. Throughout this section we assume that the hypotheses of Theorem I.1 hold.

Lemma III.2. Let $(nz_{1,n}, \dots, nz_{r,n})$ be the r largest singular values of \widehat{N} . Then, as $n \rightarrow \infty$, $z_{i,n} \rightarrow z_i$ almost surely, where, for $\Sigma_i^2 > \sigma^2/p$,

$$z_i = p \Sigma_i \left\{ \alpha \left(\frac{\sigma^2}{p \Sigma_i^2} + \frac{1}{\sqrt{\alpha}} \right) \left(\frac{\sigma^2}{p \Sigma_i^2} + \sqrt{\alpha} \right) \right\}^{1/2}, \quad (6)$$

and $z_i = \sigma \sqrt{p \alpha^{1/2} (1 + \sqrt{\alpha})}$ for $\Sigma_i^2 \leq \sigma^2/p$.

Further, let $X \in \mathbb{R}^{m \times r}$ and $Y \in \mathbb{R}^{n \times r}$ be the matrices whose columns are the first r , right and left, singular vectors

of \widehat{N} . Then there exists a sequence of $r \times r$ orthogonal matrices Q_n such that, almost surely $\|\frac{1}{\sqrt{m}} U^T X - A Q_n\|_F \rightarrow 0$, $\|\frac{1}{\sqrt{n}} V^T Y - B Q_n\|_F \rightarrow 0$ with $A = \text{diag}(a_1, \dots, a_r)$, $B = \text{diag}(b_1, \dots, b_r)$ and

$$\begin{aligned} a_i^2 &= \left(1 - \frac{\sigma^4}{p^2 \Sigma_i^4} \right) \left(1 + \frac{\sqrt{\alpha} \sigma^2}{p \Sigma_i^2} \right)^{-1}, \\ b_i^2 &= \left(1 - \frac{\sigma^4}{p^2 \Sigma_i^4} \right) \left(1 + \frac{\sigma^2}{p \sqrt{\alpha} \Sigma_i^2} \right)^{-1}, \end{aligned} \quad (7)$$

for $\Sigma_i^2 > \sigma^2/p$, while $a_i = b_i = 0$ otherwise.

Proof: Due to space limitations, we will focus here on the case $\Sigma_1, \dots, \Sigma_r > \sigma^2/p$. The general proof proceeds along the same lines, and we defer it to [4].

Notice that W^E is an $m \times n$ matrix with i.i.d. entries with variance $\sqrt{mn} \sigma^2 p$ and fourth moment bounded by $C n^2$. It is therefore sufficient to prove our claim for $p = 1$ and then rescale Σ by p and σ by \sqrt{p} . We will also assume that, without loss of generality, $m \geq n$.

Let \widehat{Z} be an $r \times r$ diagonal matrix containing the eigenvalues $(nz_{n,1}, \dots, nz_{n,r})$. The eigenvalue equations read

$$U \widehat{\beta}_y + W Y - X \widehat{Z} = 0, \quad (8)$$

$$V \widehat{\beta}_x + W^T X - Y \widehat{Z} = 0. \quad (9)$$

where we defined $\hat{\beta}_x \equiv \Sigma U^T X$, $\hat{\beta}_y \equiv \Sigma V^T Y \in \mathbb{R}^{r \times r}$. By singular value decomposition we can write $W = L \text{diag}(w_1, w_2, \dots, w_n) R^T$, with $L^T L = I_{m \times m}$, $R^T R = I_{n \times n}$.

Let $u_i^T, x_i^T, v_i^T, y_i^T \in \mathbb{R}^r$ be the i -th row of -respectively- $L^T U, L^T X, R^T V, R^T Y$. In this basis equations (8) and (9) read

$$\begin{aligned} u_i^T \hat{\beta}_y + w_i y_i^T - x_i^T \hat{Z} &= 0, & i \in [n], \\ u_i^T \hat{\beta}_y - x_i^T \hat{Z} &= 0, & i \in [m] \setminus [n], \\ v_i^T \hat{\beta}_x + w_i x_i^T - y_i^T \hat{Z} &= 0, & i \in [n]. \end{aligned}$$

These can be solved to get

$$\begin{aligned} x_i^T &= (u_i^T \hat{\beta}_y \hat{Z} + w_i v_i^T \hat{\beta}_x) (Z^2 - w_i^2)^{-1}, & i \in [n], \\ x_i^T &= u_i^T \hat{\beta}_y \hat{Z}^{-1}, & i \in [m] \setminus [n], \\ y_i^T &= (v_i^T \hat{\beta}_x \hat{Z} + w_i u_i^T \hat{\beta}_y) (\hat{Z}^2 - w_i^2)^{-1}, & i \in [n]. \end{aligned} \quad (10)$$

By definition $\Sigma^{-1} \hat{\beta}_x = \sum_{i=1}^m u_i x_i^T$, and $\Sigma^{-1} \hat{\beta}_y = \sum_{i=1}^n v_i y_i^T$, whence

$$\begin{aligned} \Sigma^{-1} \hat{\beta}_x &= \sum_{i=1}^n u_i (u_i^T \hat{\beta}_y \hat{Z} + w_i v_i^T \hat{\beta}_x) (\hat{Z}^2 - w_i^2)^{-1} \\ &\quad + \sum_{i=n+1}^m u_i u_i^T \hat{\beta}_y \hat{Z}^{-1}, \end{aligned} \quad (11)$$

$$\Sigma^{-1} \hat{\beta}_y = \sum_{i=1}^n v_i (v_i^T \hat{\beta}_x \hat{Z} + w_i u_i^T \hat{\beta}_y) (\hat{Z}^2 - w_i^2)^{-1}. \quad (12)$$

Let $\lambda = w_i^2 \alpha^{1/2} / (m^2 \sigma^2)$. Then, it is a well known fact [13] that as $n \rightarrow \infty$ the empirical law of the λ_i 's converges weakly almost surely to the Marcenko-Pastur law, with density $\rho(\lambda) = \alpha \sqrt{(\lambda - c_-^2)(c_+^2 - \lambda)} / (2\pi\lambda)$, with $c_{\pm} = 1 \pm \alpha^{-1/2}$.

Let $\beta_x = \hat{\beta}_x / \sqrt{m}$, $\beta_y = \hat{\beta}_y / \sqrt{n}$, $Z = \hat{Z} / n$. A priori, it is not clear that the sequence (β_x, β_y, Z) -dependent on n -converges. However, it is immediate to show that the sequence is tight, and hence we can restrict ourselves to a subsequence $\Xi \equiv \{n_i\}_{i \in \mathbb{N}}$ along which a limit exists. Eventually we will show that the limit does not depend on the subsequence, apart, possibly, from the rotation Q_n . Hence we shall denote the subsequential limit, by an abuse of notation, as (β_x, β_y, Z) .

Consider now a such a convergent subsequence. It is possible to show that $\Sigma_i^2 > \sigma^2/p$ implies $Z_{ii}^2 > \alpha^{3/2} \sigma^2 c_+(\alpha)^2 + \delta$ for some positive δ . Since almost surely as $n \rightarrow \infty$, $w_i^2 < \alpha^{3/2} \sigma^2 c_+(\alpha)^2 + \delta/2$ for all i , for all purposes the summands on the rhs of Eqs. (11), (12) can be replaced by uniformly continuous, bounded functions of the limiting eigenvalues λ_i . Further, each entry of u_i (resp. v_i) is just a single coordinate of the left (right) singular vectors of the random matrix W . Using Theorem 1 in [12], it follows that any subsequential limit satisfies the equations

$$\beta_x = \Sigma \beta_y \left\{ Z \int (Z^2 - \alpha^{3/2} \sigma^2 \lambda)^{-1} \rho(\lambda) d\lambda + (\alpha - 1) Z^{-1} \right\}, \quad (13)$$

$$\beta_y = \Sigma \beta_x \left\{ Z \int (Z^2 - \alpha^{3/2} \sigma^2 \lambda)^{-1} \rho(\lambda) d\lambda \right\}, \quad (14)$$

Solving for β_y , we get an equation of the form

$$\Sigma^{-2} \beta_y = \beta_y f(Z) \quad (15)$$

where $f(\cdot)$ is a function that can be given explicitly using the Stieltjis transform of the measure $\rho(\lambda) d\lambda$. Equation (15) implies that β_y is block diagonal according to the degeneracy pattern of Σ . Considering each block, either β_y vanishes in the block (a case that can be excluded using $\Sigma_i^2 > \sigma^2/p$) or $\Sigma_i^{-2} = f(Z_{ii})$ in the block. Solving for Z_{ii} shows that the eigenvalues are uniquely determined (independent of the subsequence) and given by Eq. (6).

In order to determine β_x and β_y first observe that, since $I_{r \times r} = Y^T Y = \sum_{i=1}^n y_i y_i^T$, we have, using Eq. (10)

$$\begin{aligned} I_{r \times r} &= \sum_{i=1}^n (\hat{Z}^2 - w_i^2)^{-1} (\hat{Z} \hat{\beta}_x^T v_i + w_i \hat{\beta}_y^T u_i) \\ &\quad (v_i^T \hat{\beta}_x \hat{Z} + w_i u_i^T \hat{\beta}_y) (\hat{Z}^2 - w_i^2)^{-1}. \end{aligned}$$

In the limit $n \rightarrow \infty$, and assuming a convergent subsequence for (Z, β_x, β_y) , this sum can be computed as above. After

$$\begin{aligned} I_{r \times r} &= \left\{ \int \frac{Z^2}{(Z^2 - \alpha^{3/2} \sigma^2 \lambda)^2} \rho(\lambda) d\lambda \right\} C_x \\ &\quad + \left\{ \int \frac{\alpha^{3/2} \sigma^2 \lambda}{(Z^2 - \alpha^{3/2} \sigma^2 \lambda)^2} \rho(\lambda) d\lambda \right\} C_y, \end{aligned}$$

where $C_x = \beta_x^T \beta_x$, $C_y = \beta_y^T \beta_y$ and the functions of Z on the rhs are defined as standard analytic functions of matrices.

Using Eqs. (13), (14) and solving the above, we get $C_x = \text{diag}(\Sigma_1^2 a_1^2, \dots, \Sigma_r^2 a_r^2)$, and $B_y = \text{diag}(\Sigma_1^2 b_1^2, \dots, \Sigma_r^2 b_r^2)$. We already concluded that β_x and β_y are block diagonals with blocks in correspondence with the degeneracy pattern of Σ . Since $\beta_x^T \beta_x = C_x$ and $\beta_y^T \beta_y = C_y$ are diagonal, with the same degeneracy pattern, it follows that, inside each block of size d , each of β_x and β_y is proportional to a $d \times d$ orthogonal matrix. Therefore $\beta_x = \Sigma A Q_s$, $\beta_y = \Sigma B Q'_s$, for some orthogonal matrices Q_s, Q'_s . Also, using equation (13) one can prove that $Q_s = Q'_s$.

Notice, by the above argument A, B are uniquely fixed by our construction. On the other hand Q_s might depend on the subsequence Ξ . Since our statement allows for a sequence of rotations Q_n , that depend on n , the eventual subsequence dependence of Q_s can be factored out. ■

It is useful to point out a straightforward consequence of the above.

Corollary III.3. *There exists a sequence of orthogonal matrices $Q_n \in \mathbb{R}^{r \times r}$ such that, almost surely,*

$$\lim_{n \rightarrow \infty} \left\| \frac{1}{\sqrt{mn}} X^T U \Sigma V^T Y - Q_n D Q_n^T \right\|_F = 0, \quad (16)$$

with $D = \text{diag}(\Sigma_1 a_1 b_1, \dots, \Sigma_r a_r b_r)$.

B. The effect of missing entries

The proof of Theorem I.1 is completed by establishing a relation between the singular vectors X_0, Y_0 of N^E and the singular vectors X and Y of \hat{N} .

Lemma III.4. Let $k \leq r$ be the largest integer such that $\Sigma_1 \geq \dots \geq \Sigma_k > \sigma^2/p$, and denote by $X_0^{(k)}, Y_0^{(k)}, X^{(k)}$, and $Y^{(k)}$ the matrices containing the first k columns of X_0, Y_0, X , and Y , respectively. Let $X_0^{(k)} = X^{(k)}S_x + X_\perp^{(k)}, Y_0^{(k)} = Y^{(k)}S_y + Y_\perp^{(k)}$ where $(X_\perp^{(k)})^T X^{(k)} = 0, (Y_\perp^{(k)})^T Y^{(k)} = 0$ and $S_x, S_y \in \mathbb{R}^{r \times r}$. Then there exists a numerical constant $C = C(\Sigma_i, \sigma^2, \alpha, M_{\max})$, such that, with high probability,

$$\|X_\perp^{(k)}\|_F^2, \|Y_\perp^{(k)}\|_F^2 \leq Cr \sqrt{\frac{1}{n}}, \quad (17)$$

with probability approaching 1 as $n \rightarrow \infty$.

Proof: We will prove our claim for the right singular vector Y , since the left case is completely analogous. Further we will drop the superscript k to lighten the notation.

We start by noticing that $\|N^E Y_0\|_F^2 = \sum_{a=1}^k (n\tilde{z}_{a,n})^2$, where $n\tilde{z}_{a,n}$ are the singular values of N^E . Using Lemma 3.2 in [2] which bounds $\|M^E - pM\|_2 = \|N^E - \hat{N}\|_2$, we get

$$\|N^E Y_0\|_F^2 \geq \sum_{a=1}^k (nz_{a,n} - CM_{\max}\sqrt{pn})^2. \quad (18)$$

On the other hand $\|N^E Y_0\|_F \leq \|\hat{N}Y_0\|_F + \|N^E - \hat{N}\|_2 \|Y_0\|_F$. Further by letting $S_y = L_y \Theta_y R_y^T$, for L_y, R_y orthogonal matrices, we get $\|\hat{N}Y_0\|_F^2 = \|\hat{N}Y L_y \Theta_y\|_F^2 + \|\hat{N}Y_\perp\|_F^2$. Since $Y_0^T Y_0 = I_{k \times k}$, we have $I_{k \times k} = R_y \Theta_y^T \Theta_y R_y^T + Y_\perp^T Y_\perp$, and therefore

$$\begin{aligned} \|\hat{N}Y_0\|_F^2 &= \|\hat{N}Y L_y\|_F^2 - \|\hat{N}Y L_y R_y^T Y_\perp^T\|_F^2 + \|\hat{N}Y_\perp\|_F^2 \\ &\leq n^2 \sum_{a=1}^k z_{a,n}^2 - n^2 z_{k,n}^2 \|Y_\perp\|_F^2 \\ &\quad + n^2 p \sigma^2 \alpha (c_+(\alpha) + \delta) \|Y_\perp\|_F^2 \\ &= n^2 \sum_{a=1}^k z_{a,n}^2 - n^2 e_y \|Y_\perp\|_F^2, \end{aligned}$$

where $e_y \equiv z_{k,n}^2 - p\sigma^2\alpha(c_+(\alpha) + \delta)$, and used the inequality $\|\hat{N}Y_\perp\|_F^2 \leq n^2 p \sigma^2 \alpha (c_+(\alpha) + \delta) \|Y_\perp\|_F^2$ which holds for all $\delta > 0$ asymptotically almost surely as $n \rightarrow \infty$ (by an immediate generalization of Lemma III.2). It is simple to check that $\Sigma_k \geq \sigma^2/p$ implies $e_y > 0$.

Using triangular inequality, Lemma 3.2 in [2], we get

$$\begin{aligned} \|NY_0\|_F^2 &\leq n^2 \sum_{a=1}^r z_{a,n}^2 - n^2 e_y \|Y_\perp\|_F^2 + Cn p \alpha^{3/2} M_{\max}^2 r \\ &\quad + 2Cn \sqrt{np} \alpha^{3/4} M_{\max} \sqrt{r} \|z\|, \end{aligned}$$

which, combined with equation (18), implies the thesis. \blacksquare

Proof of Theorem I.1: We now turn to upper bounding the right hand side of Eq. (5). Let k be defined as in the last lemma. Notice that by Lemma III.2, $X^T(U\Sigma V^T)Y$ is well approximated by $(X^{(k)})^T(U\Sigma V^T)Y^{(k)}$. Analogously, it can be proved that $X_0^T(U\Sigma V^T)Y_0$ is well approximated by $(X_0^{(k)})^T(U\Sigma V^T)Y_0^{(k)}$. Due to space limitations, we will omit this technical step and thus focus here on the case $k = r$

(equivalently, neglect the error incurred by this approximation).

Using Lemma III.4 to bound the contribution of X_\perp, Y_\perp , we have

$$\begin{aligned} &\langle X_0^T(U\Sigma V^T)Y_0, X_0^T N^E Y_0 \rangle \\ &= \langle S_x^T X^T(U\Sigma V^T)Y S_y, X_0^T N^E Y_0 \rangle (1 + o_n(1)) \\ &= \langle X^T(U\Sigma V^T)Y, S_x^T X_0^T N^E Y_0 S_y \rangle (1 + o_n(1)). \quad (19) \end{aligned}$$

Further $X_0^T N^E Y_0 = X_0^T \hat{N} Y_0 + X_0^T (N^E - \hat{N}) Y_0$ and, using once more the bound in Lemma 3.2 of [2], that implies $|X_0^T (N^E - \hat{N}) Y_0| \leq Cr \sqrt{nrp}$, we get

$$\begin{aligned} S_x^T X_0^T N^E Y_0 S_y &= L_x \Theta_x^2 L_x^T X^T \hat{N} Y R_y \Theta_y^2 R_y^T + E_1 \\ &= Z + E_2, \end{aligned}$$

where we recall that Z is the diagonal matrix with entries given by the singular values of \hat{N} , and $\|E_1\|_F^2, \|E_2\|_F^2 \leq C(p, r)\sqrt{n}$. Using this estimate in Eq. (19), together with the result in Lemma III.2, we finally get

$$\frac{\langle X_0^T(U\Sigma V^T)Y_0, X_0^T N^E Y_0 \rangle}{\sqrt{mn} \|X_0^T N^E Y_0\|_F} \geq \frac{\sum_{k=1}^r \Sigma_k a_k b_k z_k}{\sqrt{\alpha} \|z\|} - o_n(1),$$

which implies the thesis after simple algebraic manipulations \blacksquare

ACKNOWLEDGEMENTS

We are grateful to T. Hastie, R. Mazumder and R. Tibshirani for stimulating discussions, and for making available their data. This work was supported by a Terman fellowship, and the NSF grants CCF-0743978 and DMS-0806211.

REFERENCES

- [1] "Netflix prize," <http://www.netflixprize.com/>.
- [2] R. H. Keshavan, A. Montanari, and S. Oh, "Matrix completion from a few entries," January 2009, arxiv:0901.3150.
- [3] —, "Matrix completion from noisy entries," June 2009, arXiv:0906.2027.
- [4] R. H. Keshavan and A. Montanari, "Regularization for matrix completion," 2010, journal version, in preparation.
- [5] N. Srebro, J. Rennie, and T. Jaakkola, "Maximum margin matrix factorization," in *Advances in Neural Information Processing Systems 17*, 2005.
- [6] J. Rennie and N. Srebro, "Fast maximum margin matrix factorization for collaborative prediction," in *22nd International Conference on Machine Learning*, 2005.
- [7] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Found. of Comput. Math.*, vol. 9, no. 6, pp. 717 – 772, 2009.
- [8] E. J. Candès and Y. Plan, "Matrix completion with noise," 2009, arXiv:0903.3131.
- [9] R. Mazumder, T. Hastie, and R. Tibshirani, "Spectral regularization algorithms for learning large incomplete matrices," 2009, submitted.
- [10] I. M. Johnstone and A. Y. Lu, "On consistency and sparsity for principal component analysis in high dimension," *J. Amer. Stat. Assoc.*, vol. 104, pp. 682–693, 2009.
- [11] M. Capitaine, C. Donati-Martin, and D. Féral, "The largest eigenvalue of finite rank deformation of large wigner matrices: convergence and non-universality of the fluctuations," *Ann. Probab.*, vol. 37, pp. 1–47, 2009.
- [12] Z.D.Bai, B.Q.Miao, and G.M.Pan, "On asymptotics of eigenvectors of large sample covariance matrices," *Ann. of Probab.*, vol. 35, pp. 1532–1572, 2007.
- [13] J. Silverstein and Z. Bai, "On the empirical distribution of eigenvalues of a class of large-dimensional random matrices," *J. Multivariate Anal.*, vol. 54, pp. 175–192, 1995.