# Random Consensus Robust PCA

#### Daniel Pimentel-Alarcón & Robert Nowak UNIVERSITY of WISCONSIN-MADISON

#### Abstract

This paper presents R2PCA, a random consensus method for robust principal component analysis. R2PCA takes RANSAC's principle of using as little data as possible one step further. It iteratively selects small subsets of the data to identify *pieces* of the principal components, to then *stitch* them together. We show that if the principal components are in general position and the errors are sufficiently sparse, R2PCA will exactly recover the principal components with probability 1, in lieu of assumptions on coherence or the distribution of the sparse errors, and even under adversarial settings. R2PCA enjoys many advantages: it works well under noise, its computational complexity scales linearly in the ambient dimension, it is easily parallelizable, and due to its low sample complexity, it can be used in settings where data is so large it cannot even be stored in memory. We complement our theoretical findings with synthetic and real data experiments showing that R2PCA outperforms state-of-the-art methods in a broad range of settings.

# 1 Introduction

In many relevant applications one aims to find a lowdimensional subspace that approximates a large data matrix **M**. Principal Component Analysis (PCA) is one of the most widely used techniques for this purpose. Unfortunately, a single grossly corrupted datum can severely compromise its performance. Hence there is a wide variety of approaches to make PCA robust. Examples include M-estimators [1], random sampling [2], influence function techniques [3], alternating minimization [4] and convex relaxations [5–10]. Other approaches use convex optimization methods on subsets of the data (e.g., full rows and full columns) to improve computational complexity [11, 12].

One of the most natural and widely used algorithms for robust estimation is random sample consensus (RANSAC) [2]. RANSAC is simple yet powerful. It is popular because it makes almost no assumptions about the data, and it does not require unrealistic conditions to succeed. It has theoretical guarantees, works well in practice, and has enjoyed many improvements since its inception, e.g., [13–15]. The RANSAC version of PCA iteratively selects a few columns in **M** to define a candidate subspace, until it identifies a subspace that agrees with other columns in **M**. This will successfully identify the subspace that we are looking for, as long as **M** has enough uncorrupted columns.

But in many modern applications, such as image processing and networks data analysis, every column in  $\mathbf{M}$  may have a few grossly corrupted entries. This makes all columns in  $\mathbf{M}$  outliers, hence standard robust methods like RANSAC are no longer applicable. In this setting  $\mathbf{M}$  can be better modeled as the sum of a low-rank matrix  $\mathbf{L}$  and a sparse matrix  $\mathbf{S}$  representing the corrupted entries. The goal is to identify the subspace U spanned by the columns in  $\mathbf{L}$ . This problem is often called robust PCA (RPCA) [6]. The last decades have seen great approaches to this problem [16], yet it remained unclear how to extend RANSAC's principles to this setting [3].

The main contribution of this paper is R2PCA: a random consensus algorithm for RPCA. R2PCA takes RANSAC's principle of using as little data as possible one step further. It iteratively selects small subsets of the entries in  $\mathbf{M}$  to identify *pieces* of the subspace U. This process is repeated until we identify enough pieces to *stitch* them together and recover the whole subspace. The key idea behind R2PCA is that subspaces can be easily and efficiently recovered from a few of its canonical projections [17]. These canonical projections are the pieces that R2PCA aims to identify. See Figure 1 for some intuition.

Our main result shows that R2PCA will exactly recover the subspace that we are looking for with probabil-

Appearing in Proceedings of the 20<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2017, Fort Lauderdale, Flordia, USA. JMLR: W&CP volume 54. Copyright 2017 by the authors.



Figure 1: **First:** Each column in a rank-r matrix **L** corresponds to a point in an r-dimensional subspace U. In these figures, U is a 1-dimensional subspace (line) in general position, and the columns in **L** are drawn generically from U, that is, independently according to an absolutely continuous distribution with respect to the Lebesgue measure on U. For example, according to a gaussian distribution on U. **Second:** Adding **S** equates to corrupting some coordinates of some columns in **L**. In this figure each point is corrupted in only one coordinate. As a result, the corrupted points no longer lie in U. So, how can we identify U from these corrupted points? **Third:** The key idea behind R2PCA is that since errors are sparse, if we only focus on a few coordinates of a few columns at a time, it is likely that the selected columns are uncorrupted on the selected coordinates. We can verify whether this is the case, because the projections of the selected columns onto the selected coordinates will agree if and only if the columns. The projections of both columns onto the (x, y) coordinates agree. Namely, they both lie in  $U_{\omega_1}$ . Hence we can be sure that the (x, y) coordinates of these columns are uncorrupted, and that  $U_{\omega_1}$  is actually equal to the projection of the subspace U that we aim to identify. Last: We can repeat this procedure for different sets of coordinates and columns, until we obtain enough projections to reconstruct the whole subspace.

ity 1, as long as  $\mathbf{M}$  is generic, and the corrupted entries are sufficiently sparse. In contrast to popular convex relaxation methods (e.g., [5–12]), our results make no assumptions about coherence or the distribution of the sparse errors. In fact, our results hold even under adversarial settings where the errors are purposely located to complicate success. In its noisy variant, R2PCA can consistently estimate the desired subspace within the noise level. The computational complexity of R2PCA scales linearly in the ambient dimension. In addition, R2PCA enjoys many of RANSAC's advantages, and many of RANSAC's improvements can be easily adapted to R2PCA. For instance, R2PCA can run in parallel, with different computers searching for different pieces (canonical projections) of the subspace. This can greatly reduce computation time, which is of great interest in general, and paramount in realtime applications. Furthermore, R2PCA's principle of studying subspaces by pieces also improves computational and sample complexity. This is because R2PCA only uses small subsets of the data at a time. This is of particular importance in settings where **M** is so large it cannot even be stored in memory. We complement our theoretical findings with synthetic and real data experiments showing that R2PCA outperforms stateof-the-art methods, both in terms of speed and accuracy, in a broad range of settings.

#### 2 Model and Main Results

Suppose we observe a  $d \times n$  data matrix **M**, given by

$$\mathbf{M} = \mathbf{L} + \mathbf{S}, \tag{1}$$

where **L** is rank-*r* and **S** is sparse. The goal is to identify the *r*-dimensional subspace *U* spanned by the columns of **L**, or slightly more demanding, determine **S** and **L**. Consider the following assumptions, where  $\operatorname{Gr}(r, \mathbb{R}^d)$  denotes the Grassmannian manifold of *r*dimensional subspaces in  $\mathbb{R}^d$ , and  $\|\cdot\|_0$  denotes the  $\ell_0$ -norm, given by the number of nonzero entries.

- (A1) U is drawn according to an absolutely continuous distribution with respect to the uniform measure over  $\operatorname{Gr}(r, \mathbb{R}^d)$ .
- (A2) The columns of L are drawn independently according to an absolutely continuous distribution with respect to the Lebesgue measure on U.
- (A3) The nonzero entries in **S** are drawn independently according to an absolutely continuous distribution with respect to the Lebesgue measure on  $\mathbb{R}^{\|\mathbf{S}\|_0}$ .
- (A4) S has at most  $\frac{n-r}{2(r+1)^{\alpha}}$  nonzero entries per row and at most  $\frac{d-r}{2(r+1)^{\alpha-1}}$  nonzero entries per column, with  $\alpha \geq 1$ .

A1 requires that U is a subspace in general position, and A2 requires that the columns in L are drawn generically from this subspace. Together, A1 and A2 require that L is a generic rank-r matrix. Similarly, A3 requires that S is a generic sparse matrix. See Section 4 for a further discussion about our assumptions and their relation to other typical assumptions from the literature.

A4 requires that M has at most  $O(n/r^{\alpha})$  corrupted entries per row and at most  $O(d/r^{\alpha-1})$  corrupted entries per column. Notice that since decomposing M is the same as decomposing  $\mathbf{M}^{\mathsf{T}}$ , assumption A4 can be interchanged in terms of rows and columns. A4 is a reasonable assumption because in most problems where this setting arises, S is sparse and  $r \ll d, n$ , whence A4 allows a large number of corrupted entries in M. The parameter  $\alpha$  determines the sparsity level of S, which in turn determines the computational complexity of R2PCA. The larger  $\alpha$ , the sparser S, and the faster R2PCA will succeed.

The main result of this paper is the next theorem. It states that R2PCA (Algorithm 1 below) will exactly recover U,  $\mathbf{L}$  and  $\mathbf{S}$  on  $\mathcal{O}((d+n)2^{r^{2-\alpha}})$  iterations (linear in d and n). In the extreme case where  $\mathbf{S}$  has too many errors ( $\alpha = 1$ ), R2PCA could require exponential time in r. But if  $\mathbf{S}$  is sufficiently sparse ( $\alpha \geq 2$ ), R2PCA will succeed in linear time in r. This is true even in the adversarial setting where the errors are purposely located to complicate success. In other words, the computational complexity in Theorem 1 considers the worst-case scenario. So in practice, as shown in our experiments, R2PCA can be much faster and allow a much larger number of corrupted entries. The proof is given in Appendix A.

**Theorem 1.** Let A1-A4 hold. Let  $\hat{\mathbf{U}}$ ,  $\hat{\mathbf{L}}$  and  $\hat{\mathbf{S}}$  be the output of R2PCA. Then span{ $\hat{\mathbf{U}}$ } = U,  $\hat{\mathbf{L}} = \mathbf{L}$ and  $\hat{\mathbf{S}} = \mathbf{S}$  with probability 1. Furthermore, the expected number of iterations required by R2PCA is upper bounded by  $(d + n - r)2^{(r+1)^{2-\alpha}}$ .

**Remark 1.** Throughout the paper we assume that the rank r is known. If this is not the case, r can be estimated by iteratively selecting  $(\tau + 1) \times (\tau + 1)$  minors in  $\mathbf{M}$ , and verifying their rank (or their singular value decomposition in the noisy setting). Under A1-A4, with probability 1 there will be a  $(\tau + 1) \times (\tau + 1)$ , rank- $\tau$  minor in **M** if and only if  $r = \tau$ . So we can start with  $\tau = 1$ . If we find a  $2 \times 2$  minor in **M**, we know that r = 1. If there exists no such minor, we know that  $r \geq 2$ . We can iteratively repeat this process until we find a  $(\tau + 1) \times (\tau + 1)$  minor of rank- $\tau$ . A1-A4 imply that if  $r = \tau$ , then for every  $\boldsymbol{\omega} \subset \{1, \ldots, d\}$ with  $\tau + 1$  entries,  $\mathbf{M}_{\omega}$ , will contain a  $(\tau + 1) \times (\tau + 1)$ , rank- $\tau$  minor. Furthermore, using the same reasoning as in the proof of Theorem 1, one can show that on expectation, it would take no more than  $O(2^{r^{2-\alpha}})$  trials to find such a minor (recall that  $\alpha \geq 1$  determines the sparsity level in  $\mathbf{S}$ ).

### 3 Algorithm

In this section we introduce R2PCA in its most basic setting (Algorithm 1). In Section 5 we discuss how to generalize it to noisy settings. From a high level perspective, R2PCA searches for small subsets of uncontaminated data in  $\mathbf{M}$  to obtain *pieces* of U to then *stitch* them together. Once U is known, R2PCA searches for a few uncontaminated entries in each column of  $\mathbf{M}$  to recover the coefficients of  $\mathbf{L}$ . Once  $\mathbf{L}$  is known,  $\mathbf{S}$  can be trivially recovered through (1).

The key idea behind R2PCA is that subspaces can be exactly and efficiently recovered from a few of its canonical projections [17]. So rather than attempting to identify U directly, we will aim to identify small projections of U, knowing that there is a simple way to *stitch* these projections together to recover U. More precisely, let  $\Omega$  be a  $d \times N$  matrix with exactly r+1nonzero entries per column, and let  $\omega_i \subset \{1, 2, \ldots, d\}$ index the nonzero entries in the  $i^{\text{th}}$  column of  $\Omega$ .  $\omega_i$ indicates the canonical coordinates involved in the  $i^{\text{th}}$ projection that we will aim to identify. For any subspace, matrix or vector that is compatible with  $\boldsymbol{\omega}_i$ , we will use the subscript  $\omega_i$  to denote its restriction to the coordinates/rows in  $\omega_i$ . For example,  $\mathbf{M}_{\boldsymbol{\omega}_i} \in \mathbb{R}^{(r+1) \times n}$  denotes the restriction of **M** to the rows in  $\boldsymbol{\omega}_i$ , and  $U_{\boldsymbol{\omega}_i} \subset \mathbb{R}^{r+1}$  denotes the projection of U onto the coordinates in  $\omega_i$ .

Our goal is to identify a collection of projections  $\{U_{\boldsymbol{\omega}_i}\}_{i=1}^N$  such that U can be recovered from these projections. Whether this is the case depends on the  $\omega_i$ 's, i.e., on  $\Omega$ . Fortunately, Theorem 1 in [17] specifies the conditions on  $\mathbf{\Omega}$  to guarantee that U can be recovered from  $\{U_{\boldsymbol{\omega}_i}\}_{i=1}^N$ . To present this result, let us introduce the matrix A. A1 implies that with probability 1,  $U_{\boldsymbol{\omega}_i}$  is a hyperplane, i.e., an *r*-dimensional subspace in  $\mathbb{R}^{r+1}$ . As such, it is characterized by its orthogonal direction, which we will call  $\mathbf{a}_{\boldsymbol{\omega}_i}$ . More precisely, let  $\mathbf{a}_{\boldsymbol{\omega}_i} \in \mathbb{R}^{r+1}$  be a nonzero vector in ker  $U_{\boldsymbol{\omega}_i}$ , and let  $\mathbf{a}_i$ be the vector in  $\mathbb{R}^d$  with the entries of  $\mathbf{a}_{\omega_i}$  in the locations of  $\omega_i$ , and zeros elsewhere. Let **A** be the  $d \times N$ matrix formed with  $\{\mathbf{a}_i\}_{i=1}^N$  as columns. This way,  $\mathbf{A}$ encodes the information of the projections  $\{U_{\omega_i}\}_{i=1}^N$ . With this, we are ready to present Theorem 1 in [17], which we restate here as Lemma 1 with some adaptations to our context.

**Lemma 1** (Theorem 1 in [17]). Let **A1** hold. With probability 1,  $U = \ker \mathbf{A}^{\mathsf{T}}$  if and only if

(i) There is a matrix Ω' formed with d-r columns of
 Ω, such that every matrix formed with a subset of
 η columns in Ω' has at least η + r nonzero rows.

There exist plenty of matrices  $\Omega$  satisfying (i). For example:

$$\mathbf{\Omega} = \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{d} - r \end{bmatrix} \begin{cases} d - (r+1) \\ r+1, \end{cases}$$

where **1** and **0** denote blocks of all 1's and all 0's. One may easily verify that  $\Omega$  satisfies condition (i) by taking  $\Omega' = \Omega$ . Notice that **A** is sparse (it only has r+1 nonzero entries per column), so computing ker  $\mathbf{A}^{\mathsf{T}}$ can be done efficiently.

Lemma 1 implies that N = d - r projections are necessary and sufficient to recover U. Furthermore, it tells us which projections to look for, and how to reconstruct U from these projections. Hence, our strategy will be to select a  $d \times (d - r)$  matrix  $\Omega$  satisfying (i), then identify the projections  $\{U_{\omega_i}\}_{i=1}^{d-r}$ , and finally construct  $\mathbf{A}$  according to these projections to recover U as ker  $\mathbf{A}^{\mathsf{T}}$ .

In principle, our strategy to identify each projection is very similar to RANSAC. Given  $\omega_i$ , we iteratively select r + 1 columns of  $\mathbf{M}_{\boldsymbol{\omega}_i}$  uniformly at random. Let  $\mathbf{M}'_{\boldsymbol{\omega}_i} \in \mathbb{R}^{(r+1) \times (r+1)}$  be the matrix formed with the selected columns. span $\{\mathbf{M}'_{\boldsymbol{\omega}_i}\}$  defines a candidate projection of U onto  $\boldsymbol{\omega}_i$ . A1-A3 imply that with probability 1, span{ $\mathbf{M}'_{\boldsymbol{\omega}_i}$ } =  $U_{\boldsymbol{\omega}_i}$  if and only if  $\mathbf{M}'_{\boldsymbol{\omega}_i}$ has no corrupted entries. This will be the case if and only if rank $(\mathbf{M}'_{\boldsymbol{\omega}_i}) = r$ . We will thus verify whether  $\operatorname{rank}(\mathbf{M}'_{\omega_i}) = r$ . If this is not the case, this candidate projection will be discarded, and we will try a different  $\mathbf{M}'_{\boldsymbol{\omega}_i}$ . On the other hand, if rank $(\mathbf{M}'_{\boldsymbol{\omega}_i}) = r$ , then we know that span{ $\mathbf{M}'_{\boldsymbol{\omega}_i}$ } is the projection  $U_{\boldsymbol{\omega}_i}$  that we were looking for. In this case, we can construct  $\mathbf{a}_i$  as before. This process is repeated for each column  $\omega_i$ in  $\Omega$  to obtain A. Since  $\Omega$  satisfies (i), we know by Lemma 1 that at the end of this procedure we will have enough projections to reconstruct U as dim ker  $\mathbf{A}^{\mathsf{T}}$ .

At this point, we have already recovered U. Let  $\mathbf{U} \in \mathbb{R}^{d \times r}$  be a basis of U. We will now estimate the coefficient matrix  $\Theta \in \mathbb{R}^{r \times n}$  such that  $\mathbf{L} = \mathbf{U}\Theta$ . Let **m** be a column in **M**, and let  $\boldsymbol{\omega}$  be a subset of  $\{1, 2, \ldots, d\}$  with exactly r + 1 elements. A1-A3 imply that with probability 1,  $\mathbf{m}_{\omega}$  will lie in  $U_{\omega}$  if and only if all entries in  $\mathbf{m}_{\boldsymbol{\omega}}$  are uncorrupted. We will thus iteratively select a set  $\boldsymbol{\omega}$  indexing r+1 random entries in **m** until we find an  $\boldsymbol{\omega}$  such that  $\mathbf{m}_{\boldsymbol{\omega}} \in U_{\boldsymbol{\omega}}$ . Once we find such  $\omega$ , the coefficient vector of the corresponding column in **L** will be given by  $\boldsymbol{\theta} := (\mathbf{U}_{\omega}^{\mathsf{T}}\mathbf{U}_{\omega})^{-1}\mathbf{U}_{\omega}^{\mathsf{T}}\mathbf{m}_{\omega}$ . This process will be repeated for every column in  $\mathbf{M}$  to obtain the coefficient matrix  $\Theta$ , which together with **U** determine **L** as **U** $\Theta$ . Once **L** is known, one can trivially recover **S** through (1). R2PCA is summarized in Algorithm 1.

Algorithm 1: Random Robust PCA (R2PCA) 1 Input: Data matrix  $\mathbf{M} \in \mathbb{R}^{d \times n}$ , rank r, matrix  $\mathbf{\Omega} \in \{0, 1\}^{d \times (d-r)}$  satisfying (i).  $\mathbf{2}$ PART 1: Estimate U 3 for i = 1, 2, ..., d - r do 4  $\omega_i$  = indices of the r + 1 nonzero rows of 5 the  $i^{\text{th}}$  column in  $\Omega$ . 6 repeat 7  $\mathbf{M}'_{\boldsymbol{\omega}_i} \in \mathbb{R}^{(r+1) \times (r+1)} = r+1$  columns of 8  $\mathbf{M}_{\boldsymbol{\omega}_i}$ , selected randomly. until rank $(\mathbf{M}'_{\boldsymbol{\omega}_i}) = r.$  $\mathbf{a}_{\boldsymbol{\omega}_i} \in \mathbb{R}^{r+1} = \text{nonzero vector in ker } \mathbf{M}'_{\boldsymbol{\omega}_i}^{\mathsf{T}}.$ 9 10  $\mathbf{a}_i \in \mathbb{R}^d$  = vector with the entries of  $\mathbf{a}_{\boldsymbol{\omega}_i}$  in 11 the locations of  $\omega_i$ , and zeros 12 elsewhere. 13  $\mathbf{A} \in \mathbb{R}^{d \times (d-r)} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \mathbf{a}_{d-r}].$ 14  $\hat{\mathbf{U}} \in \mathbb{R}^{d \times r} = \text{basis of } \ker \mathbf{A}^{\mathsf{T}}.$ 1516 PART 2: Estimate  $\Theta$ for each column m in M do 17 repeat 18  $\boldsymbol{\omega} =$ subset of  $\{1, 2, \dots, d\}$  with r+119 elements, selected randomly. 20 until  $\mathbf{m}_{\boldsymbol{\omega}} \in \operatorname{span}\{\hat{\mathbf{U}}_{\boldsymbol{\omega}}\}.$ 21  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{U}}_{\boldsymbol{\omega}}^{\mathsf{T}} \hat{\mathbf{U}}_{\boldsymbol{\omega}})^{-1} \hat{\mathbf{U}}_{\boldsymbol{\omega}}^{\mathsf{T}} \mathbf{m}_{\boldsymbol{\omega}}.$ 22 Insert  $\hat{\boldsymbol{\theta}}$  into  $\hat{\boldsymbol{\Theta}}$ . 23 24 Output:  $\hat{\mathbf{U}}, \hat{\mathbf{L}} = \hat{\mathbf{U}}\hat{\Theta}, \hat{\mathbf{S}} = \mathbf{M} - \hat{\mathbf{L}}.$ 

#### 4 More about our Assumptions

Essentially, A1-A3 require that M is a combination of a generic sparse matrix, and a generic low-rank matrix. This discards pathological cases, like matrices with identical columns or exact-zero entries. Examples of these cases could arise in unnatural, cartoonlike images.

However, A1-A3 allow realistic cases, like natural images. For instance, backgrounds in natural images can be highly structured but are not perfectly constant, as there is always some degree of natural variation that is reasonably modeled by an absolutely continuous (but possibly highly inhomogeneous) distribution. For example, the sky in a natural image might be strongly biased towards blue values, but each sky pixel will have at least small variations that will make the sky not perfectly constant blue. So while these are structured images, these variations make them generic enough so that our theoretical results are applicable. This is confirmed in our real data experiments.

Furthermore, because absolutely continuous distributions may be strongly inhomogeneous, they can be used to represent highly coherent matrices (that is, matrices whose underlying subspace is highly aligned with the canonical axes). Previous theory and methods for RPCA cannot handle some of the highly coherent cases that our new theory covers and that our new algorithm handles well (as demonstrated in our experiments).

We point out that **A1-A3** do not imply coherence nor vice-versa. For example, coherence assumptions indeed allow some identical columns, or exact-zero entries. However, they rule-out cases that our theory allows. For example, consider a case where a few rows of **U** are drawn i.i.d.  $\mathcal{N}(0, \sigma_1^2)$  and many rows of **U** are drawn i.i.d.  $\mathcal{N}(0, \sigma_2^2)$ , with  $\sigma_1 \gg \sigma_2$ . This is a good model for some microscopy and astronomical applications that have a few high-intensity pixels, and many low-intensity pixels. Such **U** would yield a highly coherent matrix, which existing theory and algorithms cannot handle, while our results can (this can be confirmed in our experiments).

To sum up, our assumptions are different, not stronger nor weaker than the usual coherence assumptions, and we believe they are also more reasonable in many practical applications.

#### 5 Handling Noise

In practice, all entries in  $\mathbf{M}$  may be noisy, even the ones that are not corrupted by gross errors. We can model this as

$$\mathbf{M} = \mathbf{L} + \mathbf{S} + \mathbf{W}, \qquad (2)$$

where  $\mathbf{L}$  and  $\mathbf{S}$  are as before, and  $\mathbf{W}$  represents a noise matrix. The goal is the same as before: determine  $\mathbf{L}$  and  $\mathbf{S}$  from  $\mathbf{M}$ .

Recall that R2PCA's goal is to identify the projections  $\{U_{\boldsymbol{\omega}_i}\}_{i=1}^{d-r}$  to then reconstruct U. In the noiseless setting, we do this by iteratively selecting (r +1)  $\times$  (r+1) matrices  $\mathbf{M}'_{\boldsymbol{\omega}_i}$ , and checking their rank. If  $\operatorname{rank}(\mathbf{M}'_{\omega_i}) = r$ , then we know that all entries in  $\mathbf{M}'_{\omega_i}$ are uncorrupted, whence  $U_{\omega_i}$  is given by span $\{\mathbf{M}'_{\omega_i}\}$ . But in the noisy setting,  $\operatorname{rank}(\mathbf{M}'_{\boldsymbol{\omega}_i}) = r+1$  in general, regardless of whether these columns are corrupted by gross errors. Hence we cannot determine directly whether the columns in  $\mathbf{M}'_{\boldsymbol{\omega}_i}$  are uncorrupted by simply checking whether  $\operatorname{rank}(\mathbf{M}'_{\omega_i}) = r$ , as we did in the noiseless setting. Instead, we can check the  $(r+1)^{\text{th}}$ singular value of  $\mathbf{M}'_{\boldsymbol{\omega}_i}$ , which we will denote as  $\lambda_{r+1}$ . If  $\lambda_{r+1}$  is above the noise level, it is likely that at least one entry in  $\mathbf{M}'_{\omega_i}$  is grossly corrupted. On the other hand, if  $\lambda_{r+1}$  is within the noise level, it is likely that  $\mathbf{M}'_{\boldsymbol{\omega}_i}$  has no grossly corrupted entries, whence we can use the subspace spanned by the r leading singular

vectors of  $\mathbf{M}'_{\omega_i}$  as an estimator of  $U_{\omega_i}$ . Unfortunately, since  $\mathbf{M}'_{\omega_i}$  only has r + 1 rows and columns, the singular values and vectors of  $\mathbf{M}'_{\omega_i}$  will have a large variance. This means that  $\lambda_{r+1}$  will be above the noise level for many uncorrupted matrices  $\mathbf{M}'_{\omega_i}$ , and below the noise level for many corrupted ones. As a result, we could miss good estimators and use bad ones. Furthermore, even if  $\mathbf{M}'_{\omega_i}$  is uncorrupted, the subspace spanned by its r leading singular vectors could be far from  $U_{\omega_i}$ . As a result, our estimate of U could be very inaccurate.

But this can be improved if we use a few more entries in M so that the noise cancels out. To this end, let  $\kappa_i$ be a subset of  $\{1, 2, \ldots, d\}$  with k > r elements containing  $\boldsymbol{\omega}_i$ , and let  $\mathbf{M}'_{\boldsymbol{\kappa}_i} \in \mathbb{R}^{k \times k}$  be a matrix formed with k columns of  $\mathbf{M}_{\boldsymbol{\kappa}_i}$ . Define  $\mathbf{V}_{\boldsymbol{\kappa}_i} \in \mathbb{R}^{k \times r}$  as the matrix formed with the r leading left singular vectors of  $\mathbf{M}'_{\kappa}$ . Under mild, typical assumptions (e.g., finite second and fourth moments), if  $\mathbf{M}_{\kappa_i}$  is uncorrupted, as k grows, the  $(r+1)^{\text{th}}$  singular value of  $\mathbf{M}'_{\kappa_i}$  converges to the noise level, and span{ $\mathbf{V}_{\kappa_i}$ } converges to  $U_{\kappa_i}$ . In other words, the larger k, the better estimates of U we will obtain. On the other hand, as k grows, it is more likely that at least one entry in  $\mathbf{M}'_{\kappa_i}$  is grossly corrupted (because  $\mathbf{M}'_{\boldsymbol{\kappa}_i}$  will have more entries, each of which may be grossly corrupted), contrary to what we want. We thus want k to be large enough such that  $\mathbf{M}'_{\kappa_i}$  can be used to accurately estimate  $U_{\kappa_i}$ , but not so large that there are no matrices  $\mathbf{M}'_{\kappa_i}$  with uncorrupted entries. The fraction of corrupted entries in M determines how large k can be. Figure 3 in Section 6 shows the feasible range of k as a function of the fraction of corrupted entries in **M**.

Since  $\mathbf{M}'_{\kappa_i}$  has k > r rows, if  $\mathbf{M}'_{\kappa_i}$  is believed to be uncorrupted, we can use it to estimate several projections of U (as many as k-r). To see this, let  $\boldsymbol{v}_i$  be a subset of  $\boldsymbol{\omega}_i$  with exactly r elements. Let  $j \in \boldsymbol{\kappa}_i \setminus \boldsymbol{v}_i$ , and let  $\boldsymbol{\omega}_{ij} := \boldsymbol{v}_i \cup j$ . Observe that,  $\mathbf{V}_{\boldsymbol{\omega}_{ij}} \in \mathbb{R}^{(r+1) \times r}$  gives an estimate of  $U_{\omega_{ij}}$  through span  $\{V_{\omega_{ij}}\}$ . As before, we will store this information in the matrix  $\mathbf{A}$ . More precisely, for each  $j \in \kappa_i \setminus v_i$ , we will take a nonzero vector  $\mathbf{a}_{\omega_{ij}} \in \ker \mathbf{V}_{\omega_{ij}}^{\mathsf{T}}$ , and we will construct the vector  $\mathbf{a}_{ij} \in \mathbb{R}^d$  with the entries of  $\mathbf{a}_{\boldsymbol{\omega}_{ij}}$  in the locations of  $\omega_{ij}$ . This time, **A** will be the matrix with the  $\mathbf{a}_{ij}$ 's as columns. Since  $\omega_i = \omega_{ij}$  for some j, Lemma 1 suggests that the projections encoded in **A** should be enough to reconstruct **U**. We can thus use the matrix  $\hat{\mathbf{U}} \in \mathbb{R}^{d \times r}$  formed with the last r left singular vectors of **A** (which approximates ker  $\mathbf{A}^{\mathsf{T}}$ ) to estimate of U.

Similarly, in the second part of R2PCA we can estimate the coefficients of **L** using k entries of each column in **M**. More precisely, for each column **m** in **M**, we can iteratively select a set  $\kappa$  indexing k random entries in **m** until we find a  $\kappa$  such that  $\mathbf{m}_{\kappa}$  is close to span{ $\hat{\mathbf{U}}_{\kappa}$ } (within the noise level). If this is the case, it is likely that the entries in  $\mathbf{m}_{\kappa}$  are uncorrupted, and that  $\hat{\boldsymbol{\theta}} = (\hat{\mathbf{U}}_{\kappa}^{\mathsf{T}} \hat{\mathbf{U}}_{\kappa})^{-1} \hat{\mathbf{U}}_{\kappa}^{\mathsf{T}} \mathbf{m}_{\kappa}$  is a good estimate of the coefficient we are looking for. We repeat this process to obtain an estimate  $\hat{\boldsymbol{\Theta}}$  of  $\boldsymbol{\Theta}$ . Finally, our estimate of  $\mathbf{L}$  is given by  $\hat{\mathbf{U}}\hat{\boldsymbol{\Theta}}$ , which in turn gives an estimate of  $\mathbf{S}$  through (1). The noisy variant of R2PCA is summarized in Algorithm 2 in Appendix B.

### 6 Experiments

In this section we present a series of experiments to study the performance of R2PCA and compare it with the augmented Lagrange multiplier method for robust PCA(RPCA-ALM) [18, 19]. We found, consistent with previous reports [9, 16, 20], that the RPCA-ALM algorithm typically performed as well or better than several others (e.g., singular value thresholding [7], alternating direction method [21], accelerated proximal gradient [22] and dual method [22]).

Synthetic Data. We will first use simulations to study the performance of R2PCA in the noiseless setting, as a function of the percentage of grossly corrupted entries per row p, and the coherence of  $\mathbf{L}$ , defined as  $\mu := \frac{d}{r} \max_{1 \le i \le d} \|\mathbf{P}_U \mathbf{e}_i\|_2^2$ , where  $\mathbf{P}_U$  denotes the projection operator onto U, and  $\mathbf{e}_i$  the  $i^{\text{th}}$ canonical vector in  $\mathbb{R}^d$ . Intuitively,  $\mu$  parametrizes how aligned is U with respect to the canonical axes. In all our experiments,  $\mathbf{L}$  was a  $d \times n$ , rank-r matrix, with d = n = 100 and r = 5.

In our simulations, we first generated a  $d \times r$  random matrix **U** with  $\mathcal{N}(0,1)$  i.i.d. entries to use as basis of U. To obtain matrices with a specific coherence parameter, we simply increased the magnitude of a few entries in **U**, until it had the desired coherence. We then generated an  $r \times (r+1)(d-r)$  random matrix  $\Theta$ , also with  $\mathcal{N}(0,1)$  i.i.d. entries, to use as coefficient vectors. With this, we constructed  $\mathbf{L} = \mathbf{U}\boldsymbol{\Theta}$ . Next, we generated a  $d \times r$  matrix **S** with p percent of nonzero entries per row, selected uniformly at random. The nonzero entries in **S** are i.i.d.  $\mathcal{N}(0, 10)$ . Finally, we obtained **M** as in (1). We repeated this experiment 100times for each pair  $(p, \mu)$ , and recorded the fraction of trials that L was exactly recovered. We declared a success if the normalized error (using Frobenius norm) was below  $10^{-10}$  after at most  $10^3$  seconds. The results are summarized in Figure 2. As predicted by our theory, R2PCA performs perfectly as long as  $\mathbf{S}$  is sufficiently sparse, regardless of coherence. In contrast, other approaches rely on low coherence (e.g., [5–12]), and one can see that their performance quickly decays as coherence increases. On the other hand, as p grows, and **S** becomes less sparse, the likelihood of finding uncorrupted blocks in M quickly decays. In turn, it takes



Figure 2: Transition diagrams of the success rate (top row) and time (bottom row) for exact recovery of **L** as a function of the percentage of grossly corrupted entries per row p, and the coherence parameter  $\mu \in [1, d/r]$ . The color of each  $(p, \mu)$  pixel indicates the average over 100 trials (the lighter the better). Notice that as p grows, so does the time required to find projections, up to the point where R2PCA is unable to find enough projections to reconstruct U. Theorem 1 shows that if **M** has at most p = 7.9% corrupted entries per row (dashed line), then R2PCA can exactly recover **L**. We point out that our results hold regardless of coherence, as opposed to other algorithms, whose performance quickly decays as coherence increases.

more time to identify projections of U, up to the point where R2PCA is unable to identify enough projections to reconstruct U. Our astronomy and real-data experiments below illustrate the importance of coherent matrices and non uniformly distributed errors.

In Section 5 we presented a noisy variant of R2PCA that iteratively selects k rows of k columns of **M** to estimate U and  $\Theta$ . Its performance depends on the choice of k. If k is too small, our estimates could be very inaccurate. If k is too large, **M** may not contain enough  $k \times k$  uncorrupted blocks to obtain an estimate. The feasible range of k depends on the percentage of corrupted entries p. In our next experiment we study the performance of the noisy variant of R2PCA as a function of p and k. To obtain a noisy **M** according to (2), we generated matrices **L** and **S** as described before, and then added a  $d \times n$  random matrix **W** with  $\mathcal{N}(0, \sigma^2)$  i.i.d. entries. To measure accuracy we recorded the error of the estimated L after at most  $10^3$  seconds (using normalized Frobenius norm). We repeated this experiment 100 times for each pair (p, k)with  $\sigma = 10^{-3}$  fixed. The results, summarized in Figure 3, show the feasible range of k.

In our next simulation, we selected k = 2r, known from our previous experiment to produce reasonable results for a wide range of p, and used it to test the perfor-



Figure 3: Transition diagram of the estimation error of **L** (using the noisy variant of R2PCA) as a function of the percentage of grossly corrupted entries per row p and the parameter k, with noise level  $\sigma = 10^{-3}$ . The color of each (p, k) pixel indicates the average over 100 trials (the lighter the better). If k is too small, our estimate could be very inaccurate. If k is too large, it is less likely to find  $k \times k$  uncorrupted matrices to obtain an estimate. This figure shows the feasible range of k (white region, where R2PCA can recover **L** within the noise level), which depends on the percentage of corrupted entries p.

mance of R2PCA as a function of noise and coherence, with fixed p = 5%. We repeated this experiment 100 times for each pair  $(\sigma, \mu)$ . The results, summarized in Figure 4, show that R2PCA can consistently estimate **L** within the noise level, as long as **S** is sufficiently sparse, regardless of coherence (as opposed to other algorithms).

Astronomy and Correlated Errors. In a video, the background can be modeled as approximately lowrank, and the foreground objects (like cars or people) can be modeled as sparse errors (as they typically take only a fraction of each frame). So the sparse plus lowrank model is a natural fit for this problem. Here **M** is the matrix containing the vectorized images in the video, and the goal is to decompose it into the sparse foreground **S** and the low-rank background **L**. We now



Figure 4: Transition diagram of the estimation error of **L** as a function of the noise level  $\sigma$  and the coherence parameter  $\mu$ , with p = 5% grossly corrupted entries. The color of each  $(\sigma, \mu)$  pixel indicates the average error over 100 trials (the lighter the better). This shows that R2PCA can consistently estimate **L** within the noise level, as long as **S** is sufficiently sparse, regardless of coherence. Other algorithms can also estimate **L** within the noise level, but only for a restricted range of matrices with bounded coherence.



Figure 5: Left: One frame of a simulation of an astronomical video, composed of a background formed with  $\nu$  twinkling stars and  $\rho$  moving objects. Each object (block) moves in a random direction at a random speed over the video. **Right:** Each frame is vectorized to form the matrix **M**, which is shown negated and transposed for display purposes (i.e., we see  $1 - \mathbf{M}^{\mathsf{T}}$ ). The vertical lines correspond to the pixels of the stars. All other points correspond to the moving objects. These points are highly correlated, as is the location of an object in consecutive frames.

present an experiment where highly coherent matrices and highly correlated sparse errors arise in a very natural way: background segmentation of astronomy videos.

In this experiment we simulated astronomy videos with a background with  $\nu$  twinkling stars and  $\rho$  moving objects (see Figure 5). To this end we first generated a  $d \times r$  matrix **U**, and an  $r \times N$  matrix  $\Theta$ . with  $d = 90 \cdot 120 = 10800$ , r = 5 and N = 100. We selected  $\nu$  rows in **U** uniformly at random, and populated them with the absolute values of i.i.d.  $\mathcal{N}(0, 100)$ random variables. These entries represent the twinkling stars. All other entries in **U** were populated with the absolute values of i.d.d.  $\mathcal{N}(0,1)$  random variables. Similarly, we populated  $\Theta$  with the absolute value of i.d.d.  $\mathcal{N}(0,1)$  random variables. Next we constructed  $\mathbf{L} = \mathbf{U}\boldsymbol{\Theta}$ , and bounded its entries by 1 (i.e., we divided **L** by its maximum value). Each column of **L** represents the vectorized background of a  $90 \times 120$  frame of a video.

For each of the  $\rho$  moving objects, we selected uniformly at random: one starting point on the edge of a 90×120 frame, one starting time between {1, 2, ..., 100}, one direction, and one speed ranging from 1 to 5 pixels per frame. With this information, we created  $\rho$  objects moving across a dark background over 100 frames. Each moving object consisted of an  $r \times r$  block with  $\mathcal{N}(0, 1)$  entries. We vectorized the frames to obtain a 10800 × 100 matrix, whose entries we normalized between 0 and 1 to obtain **S**. Finally, we replaced the zero entries in **S** with the corresponding entries in **L** to obtain **M**. This way, all entries in **M** are between 0 and 1, such that the brightest star shines at a maximum intensity of 1, and so does the brightest pixel of all moving objects.

We repeated this experiment 100 times for each pair



Figure 6: Transition diagrams of the success rate (top row) and time (bottom row) for exact recovery of **L** as a function of the number of moving objects  $\rho$  and the number of twinkling stars  $\nu$  for the experiment described in Figure 5. The larger  $\rho$ , the larger proportion of corrupted entries p. The larger  $\nu$ , the lower coherence  $\mu$ . The color of each  $(\rho, \nu)$  pixel indicates the average over 100 trials (the lighter the better). The results are consistent with the experiments in Figure 2.

 $(\nu, \rho)$ . The results, summarized in Figure 6, show that R2PCA has almost perfect performance handling highly coherent matrices and highly correlated errors (as opposed to other algorithms).

**Real Data:** Microscopy and Surveillance. Finally, we evaluate the background segmentation performance of R2PCA on real data. To this end we used several microscopy videos from the Internet [25], and videos from two widely used datasets: the Wallflower dataset [23] and the I2R dataset [24]. Figure 7 shows

several examples, with more in Appendix C.

We point out that many cases of the Wallflower and the I2R datasets have low coherence. In these cases, the performance of R2PCA and RPCA-ALM is very similar. Consistent with our theory, the advantage of R2PCA becomes more evident in highly coherent cases, like our microscopy and astronomy experiments.

**Remark 2.** Notice that in all of our background experiments, R2PCA can handle a much larger fraction of gross errors than the allowed by Theorem 1. This is because Theorem 1 holds even under the worst-case scenario where the errors are purposely located to complicate success. In many applications, as in background segmentation, errors are often grouped, which tends to leave more uncorrupted blocks. This facilitates R2PCA's success.

# 7 Conclusions

In this paper we present R2PCA, a novel algorithm for robust PCA. We show that under reasonable assumptions, R2PCA will succeed with probability 1 in linear time, in lieu of assumptions on coherence or the distribution of the sparse errors. The algorithm is parallelizable and can be used in large scale settings where the dataset is too large to even store in memory. Our experiments show that R2PCA consistently outperforms state-of-the-art methods both in terms of speed and accuracy in a broad range of settings, particularly on high coherence cases.



Figure 7: Sparse (foreground) plus low-rank (background) decomposition of several microscopy videos [25] using R2PCA and RPCA-ALM [18, 19]. Notice that the background obtained by RPCA-ALM contains foreground objects, while the background obtained by R2PCA is much cleaner. This is because it these videos the background is mostly dark with a few bright regions (which implies a highly coherent subspace) and the location of the errors is highly correlated (the location of an object in consecutive frames is very similar). In contrast to other optimization methods [5–12, 18, 19], we make no assumptions about coherence or the distribution of the sparse errors, and so this does not affect our results.

### References

- [1] R. Maronna, Robust M-estimators of multivariate location and scatter, The Annals of Statistics, 1976.
- [2] M. Fischler and R. Bolles, Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography, Communications of the ACM, 1981.
- [3] F. De La Torre and M. Black, A framework for robust subspace learning, International Journal of Computer Vision, 2003.
- [4] Q. Ke and T. Kanade, Robust  $L_1$  norm factorization in the presence of outliers and missing data by alternative convex programming, IEEE Conference on Computer Vision and Pattern Recognition, 2005.
- [5] E. Candès and J. Romberg, *Sparsity and incoherence in compressive sampling*, Inverse Problems, 2007.
- [6] J. Wright, A. Ganesh, S. Rao, Y. Peng and Y. Ma, Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization, Advances in Neural Information Processing Systems, 2009.
- [7] J. Cai, E. Candès and Z. Shen, A singular value thresholding algorithm for matrix completion, SIAM Journal on Optimization, 2010.
- [8] H. Xu, C. Caramanis and S. Sanghavi, *Robust PCA via outlier pursuit*, Advances in Neural Information Processing Systems, 2010.
- [9] E. Candès, X. Li, Y. Ma and J. Wright, Robust principal component analysis?, Journal of the ACM, 2011.
- [10] V. Chandrasekaran, S. Sanghavi, P. Parrilo and A. Willsky, *Rank-sparsity incoherence for matrix decomposition*, SIAM Journal on Optimization, 2011.
- [11] L. Mackey, A. Talwalkar and M. Jordan, *Divide-and-conquer matrix factorization*, Advances in Neural Information Processing Systems, 2011.
- [12] M. Rahmani and G. Atia, A subspace learning approach for high dimensional matrix decomposition with efficient column/row sampling, International Conference on Machine Learning, 2016.
- [13] O. Chum, J. Matas and J. Kittler, *Locally optimized* RANSAC, Pattern Recognition, 2003.
- [14] O. Chum and J. Matas, *Optimal randomized* RANSAC, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2008.
- [15] A. Hast, J. Nysjö and A. Marchetti, Optimal RANSAC -Towards a repeatable algorithm for finding the optimal set, Journal of WSCG, 2013.
- [16] T. Bouwmans and E. Zahzah, Robust PCA via principal component pursuit: a review for a comparative evaluation in video surveillance, Computer Vision and Image Understanding, 2014.
- [17] D. Pimentel-Alarcón, N. Boston and R. Nowak, Deterministic conditions for subspace identifiability from incomplete sampling, IEEE International Symposium on Information Theory, 2015.
- [18] Z. Lin, M. Chen, L. Wu, and Y. Ma, *The augmented Lagrange multiplier method for exact recovery* of corrupted low-rank matrices, University of Illinois at Urbana-Champaign Technical Report, 2009.

- [19] Z. Lin, R. Liu and Z. Su, *Linearized alternating direc*tion method with adaptive penalty for low rank representation, Advances in Neural Information Processing Systems, 2011.
- [20] Y. Ma, Low-rank matrix recovery and completion via convex optimization, available at http://perception. csl.illinois.edu/matrix-rank/home.html.
- [21] X. Yuan and J. Yang, Sparse and low-rank matrix decomposition via alternating direction methods, available at http://www.optimization-online.org/ DB\_HTML/2009/11/2447.html, 2009.
- [22] Z. Lin, A. Ganesh, J. Wright, L. Wu, M. Chen and Y. Ma, Fast convex optimization algorithms for exact recovery of a corrupted low-rank matrix, Computational Advances in Multi-Sensor Adaptive Processing, 2009.
- [23] K. Toyama, J. Krumm, B. Brumitt and B. Meyers, Wallflower: principles and practice of background maintenance, IEEE International Conference on Computer Vision, 1999. Dataset available at: http://research.microsoft.com/en-us/um/ people/jckrumm/wallflower/testimages.htm
- [24] L. Li, W. Huang, I. Gu and Q. Tian, Statistical modeling of complex backgrounds for foreground object detection, IEEE Transactions on Image Processing, 2004. Dataset available at: http://perception. i2r.a-star.edu.sg/bk\_model/bk\_index.html
- [25] Microscopy videos downloaded from the (YouTube), Internet available at: https: //www.youtube.com/watch?v=iTswesT8zeo, https://www.youtube.com/watch?v=-sxrIvbqBGg, https://www.youtube.com/watch?v=sX6PTD4t-Xs, https://www.youtube.com/watch?v=Ooce\_x-SoJA, https://www.youtube.com/watch?v=R99T08FjTkc, https://www.youtube.com/watch?v=6wNd-sHU7r8, https://www.youtube.com/watch?v=gbrZhYzk480, https://www.youtube.com/watch?v=1xlNgeBc3ek.

#### Random Consensus Robust PCA Supplementary Material

#### Daniel Pimentel-Alarcón & Robert Nowak UNIVERSITY of WISCONSIN-MADISON

# A Proof of Theorem 1

In this section we give the proof of Theorem 1. Recall that  $\Omega$  is a  $d \times (d-r)$  matrix, and that  $\omega_i \subset \{1, \ldots, d\}$  indexes the r + 1 nonzero entries in the  $i^{\text{th}}$  column of  $\Omega$ . Each  $\omega_i$  indicates the coordinates of a projection of U that we aim to identify. Since R2PCA selects a matrix  $\Omega$  satisfying (i), Lemma 1 implies that if we find the projections of U onto the coordinates indicated in  $\Omega$ , then we can reconstruct U from these projections. Hence we need to show that under the assumptions of Theorem 1, R2PCA can find the projections of U onto the  $\omega_i$ 's in  $\Omega$ . To this end, we will show that R2PCA can potentially find the projection onto  $any \omega \subset \{1, 2, \ldots, d\}$  with exactly r + 1 elements.

So let  $\boldsymbol{\omega}$  be given. As discussed in Section 2, finding the projection  $U_{\boldsymbol{\omega}}$  equates to finding r+1 uncorrupted columns in  $\mathbf{M}_{\boldsymbol{\omega}}$ . So R2PCA can potentially find  $U_{\boldsymbol{\omega}}$  as long as there are r+1 uncorrupted columns in  $\mathbf{M}_{\boldsymbol{\omega}}$ . Since  $\mathbf{M}_{\boldsymbol{\omega}}$  only contains r+1 rows, **A4** implies that there are at most  $\frac{(n-r)(r+1)}{2(r+1)^{\alpha}}$  corrupted entries in  $\mathbf{M}_{\boldsymbol{\omega}}$ . In the worst-case scenario, each of these corrupted entries is located in a different column. It follows that  $\mathbf{M}_{\boldsymbol{\omega}}$  has at most  $\frac{n-r}{2(r+1)^{\alpha-1}}$  corrupted columns. Then

$$\mathsf{P}(i^{ ext{th}} ext{ column in } \mathbf{M}'_{\boldsymbol{\omega}} ext{ is uncorrupted})$$
  
  $\geq 1 - \frac{1}{2(r+1)^{\alpha-1}},$ 

which corresponds to the case where the first r columns in  $\mathbf{M}'_{\boldsymbol{\omega}}$  are uncorrupted, whence the ratio of uncorrupted columns  $((n-r) - \frac{n-r}{2(r+1)^{\alpha-1}})$  versus total remaining columns (n-r) is smallest. It follows that  $\mathsf{P}(\text{all columns in } \mathbf{M}'_{\boldsymbol{\omega}} \text{ are uncorrupted})$ 

$$\geq \left(1 - \frac{1}{2(r+1)^{\alpha-1}}\right)^{r+1}$$

$$= \left(1 - \frac{1/2}{(r+1)^{\alpha-1}}\right)^{(r+1)^{1+(\alpha-1)-(\alpha-1)}}$$

$$= \left(1 - \frac{1/2}{(r+1)^{\alpha-1}}\right)^{(r+1)^{(\alpha-1)}(r+1)^{2-\alpha}}$$

$$= \left(\left(1 - \frac{1/2}{(r+1)^{\alpha-1}}\right)^{(r+1)^{(\alpha-1)}}\right)^{(r+1)^{2-\alpha}}$$

$$\geq (1/2)^{(r+1)^{2-\alpha}}.$$

$$(1)$$

This implies that on expectation, R2PCA will require at most  $2^{(r+1)^{2-\alpha}}$  iterations to find a set of r+1 uncorrupted columns in  $\mathbf{M}_{\boldsymbol{\omega}}$ . This is true for every  $\boldsymbol{\omega}$ . Since R2PCA only searches over the  $\boldsymbol{\omega}_i$ 's in  $\boldsymbol{\Omega}$ , and since  $\boldsymbol{\Omega}$ has exactly d-r columns, it follows that on expectation, R2PCA will require at most  $(d-r)2^{(r+1)^{2-\alpha}}$  iterations to find the projections of U onto the canonical coordinates indicated by  $\boldsymbol{\Omega}$ . Since  $\boldsymbol{\Omega}$  satisfies condition (i), we know by Lemma 1 that U is given by ker  $\mathbf{A}^{\mathsf{T}}$ .

Now that U is known, let us show that R2PCA can recover **L**. Let **U** be an arbitrary basis of U. We will show that R2PCA can determine the matrix  $\Theta$  containing the coefficients of **L** in this basis, such that in the end, **L** will be given by  $\mathbf{U}\Theta$ . To this end, let **m** be a column in **M**. Observe that R2PCA can potentially find the coefficients of the corresponding column of **L** as long as there is a set  $\boldsymbol{\omega} \subset \{1, 2, \ldots, d\}$  with r+1 elements such that  $\mathbf{m}_{\boldsymbol{\omega}} \in U_{\boldsymbol{\omega}}$ . A1-A3 imply that with probability 1, this will be the case if and only there are at least r + 1 uncorrupted entries in **m**. By A4, there are at most  $\frac{d-r}{2(r+1)^{\alpha-1}}$  corrupted entries in **m**. It follows that

$$\mathsf{P}(i^{\mathrm{th}} \text{ entry in } \mathbf{m}_{\boldsymbol{\omega}} \text{ is uncorrupted})$$
  
  $\geq 1 - \frac{1}{2(r+1)^{\alpha-1}},$ 

which corresponds to the case where the first r entries

in  $\mathbf{m}_{\boldsymbol{\omega}}$  are uncorrupted, whence the ratio of uncorrupted entries  $((d-r) - \frac{d-r}{2(r+1)^{\alpha-1}})$  versus total remaining entries (d-r) is smallest. It follows that

 $P(\text{all entries in } \mathbf{m}_{\boldsymbol{\omega}} \text{ are uncorrupted})$ 

$$\geq \left(1 - \frac{1}{2(r+1)^{\alpha-1}}\right)^{r+1} \ge (1/2)^{(r+1)^{2-\alpha}},$$

where the last inequality follows by the same arithmetic manipulations as in (1). This implies that on expectation, R2PCA will require at most  $2^{(r+1)^{2-\alpha}}$  iterations to find a set of r + 1 uncorrupted entries in **m**. This is true for every **m**. Since **M** has n columns, it follows that on expectation, R2PCA will require at most  $n2^{(r+1)^{2-\alpha}}$  iterations to recover **L**. Once **L** is known, **S** can be trivially recovered as  $\mathbf{S} = \mathbf{M} - \mathbf{L}$ . This shows that on expectation, R2PCA will require at most  $(d + n - r)2^{(r+1)^{2-\alpha}}$  iterations to recover U, **L** and **S** from **M**.

# **B** Noisy Variant

In Section 5 we described a noisy variant of R2PCA. This variant iteratively selects matrices  $\mathbf{M}'_{\kappa} \in \mathbb{R}^{k \times k}$ formed with k rows of k columns of  $\mathbf{M}$ , and verifies the  $(r+1)^{\text{th}}$  singular value of  $\mathbf{M}'_{\kappa}$ . If this singular value is within the noise level, Algorithm 1 will consider  $\mathbf{M}'_{\kappa}$ uncorrupted, and use it to estimate projections of U. Otherwise Algorithm 1 will discard  $\mathbf{M}'_{\kappa}$  and keep looking. This process is repeated until there are enough projections to recover U. Once U is estimated, Algorithm 1 proceeds to estimate the coefficients of L using k entries per column of **M**. If these entries agree with the estimated subspace U, they will be considered uncorrupted, and used to estimate the coefficient of the corresponding column of L. Otherwise, Algorithm 1 will discard these entries, and select an other k. This process is repeated until we recover all the coefficients of L. This noisy variant of R2PCA is summarized in Algorithm 1.

| -  | Algorithm 1: Random Robust PCA  |  |  |
|----|---|--|--|
| -  | (R2PCA, noisy variant)  |  |  |
|    | 1 Input: Data $\mathbf{M} \in \mathbb{R}^{d \times n}$ , rank $r$ ,                 |  |  |
|    | 2   | matrix $\mathbf{\Omega} \in \{0,1\}^{d \times (d-r)}$ satisfying   |  |
|    | (   | i),  |  |
|    | 3 parameter $\kappa \in \mathbb{N}$ .   |  |  |
|    | 4 I<br>5  | for $i = 1, 2$ $d = r$ do  |  |
|    | 6   | $\omega_i = \text{indices of the } r+1 \text{ nonzero}$  |  |
|    | -   | rows of  |  |
|    | 7   | the $i^{\rm th}$ column in $\Omega$ .  |  |
|    | 8   | $\kappa_i = \text{subset of } \{1, \ldots, d\} \text{ containing}$   |  |
|    |   | $\omega_i$   |  |
|    | 9   | and $k - r + 1$ other rows   |  |
|    |   | selected randomly. repeat  |  |
|    | 10  | $\mathbf{M}_{\kappa_i} \in \mathbb{R}^{n \times n} = k \text{ columns of } \mathbf{M}_{\kappa_i},$   |  |
|    | 11  | selected randomly.   |  |
| 19 | 12  | is within the noise level  |  |
| 10 | 14  | $\mathbf{V}_{m} \in \mathbb{R}^{k \times r} = r \text{ leading singular}$  |  |
|    | 14  | $\mathbf{v}_{\kappa_i} \subset \mathbf{k}$ $-\mathbf{v}$ reading singular vectors  |  |
|    | 15  | of $\mathbf{M}'_{\kappa}$ .  |  |
|    | 16  | $\boldsymbol{v}_i = \text{subset of } \boldsymbol{\kappa}_i \text{ with exactly } r$   |  |
|    |   | elements,  |  |
|    | 17  | selected randomly.   |  |
|    | 18  | for each $j \in \kappa_i \setminus v_i$ do   |  |
|    | 19  | $\omega_{ij} := v_i \cup j.$   |  |
|    | 20  | $\mathbf{a}_{\boldsymbol{\omega}_{ij}} \in \mathbb{R}^{r+1} = \text{nonzero vector}$   |  |
|    | 21  | $\lim \ker \mathbf{V}_{\omega_{ij}}.$  |  |
|    | 22  | $\mathbf{a}_{ij} \in \mathbb{R}^{d} = \text{vector with } \mathbf{a}_{\omega_{ij}} \text{ in the locations of } \ldots$  |  |
|    | 23  | $\omega_{ij}$ , and $\omega_{ij}$  |  |
|    | 24  | elsewhere.   |  |
|    | 25  | Insert $\mathbf{a}_{ii}$ into $\mathbf{A}$ .   |  |
|    |   |  |  |
|    | 26 $\bigcup \mathbf{U} \in \mathbb{K}^{u \wedge r}$ = basis of ker A <sup>+</sup> . |  |  |
|    | 27 PART 2: Estimate $\Theta$  |  |  |
|    | 28  | for each column $\mathbf{m}$ in $\mathbf{M}$ do  |  |
|    | 29  | repeat   |  |
|    | 30  | $\kappa = $ subset of $\{1, \dots, a\}$ with $\kappa$  |  |
|    | 31  | $\mathbf{\hat{H}}$   |  |
| 99 | 32  | (within the noise level)   |  |
| 50 | 9.4   | $\hat{\boldsymbol{\mu}} = (\hat{\mathbf{I}}^{T}\hat{\mathbf{I}}^{T})^{-1}\hat{\mathbf{I}}^{T}\mathbf{m}$   |  |
|    | 34<br>97  | $\begin{bmatrix} \mathbf{U} - (\mathbf{U}_{\kappa} \mathbf{U}_{\kappa}) & \mathbf{U}_{\kappa} \mathbf{I} \\ \mathbf{I}_{nort} \hat{\boldsymbol{A}}_{nort} & \hat{\boldsymbol{A}}_{nort} & \hat{\boldsymbol{A}}_{nort} \end{bmatrix}$ |  |
|    | 35  |  |  |
|    | 36 (  | $\textbf{Dutput: } \mathbf{\hat{U}},  \mathbf{\hat{L}} = \mathbf{\hat{U}}\mathbf{\hat{\Theta}},  \mathbf{\hat{S}} = \mathbf{M} - \mathbf{\hat{L}}.$  |  |

# C Additional Results

**Microscopy Segmentation** In Section 6 we gave three examples of the background segmentation that we obtained for three microscopy videos from the Internet. Figure 1 shows more results.

Wallflower and I2R Datasets. To complement the real data experiments in Section 6, we also ran R2PCA and RPCA-ALM on the Wallflower [23] and the I2R [24] datasets. The results are summarized in Figure 2. We point out that many cases of the Wallflower and the I2R datasets have low coherence. In these cases, the performance of R2PCA and RPCA-ALM is very similar. Consistent with our theory, the advantage of R2PCA becomes more evident in highly coherent cases, like our microscopy and astronomy experiments.



Figure 1: Sparse (foreground) plus low-rank (background) decomposition of some video frames from several microscopy videos from the Internet [25] using R2PCA and RPCA-ALM [18,19]. Notice that the background obtained by RPCA-ALM contains foreground objects, while the background obtained by R2PCA is much cleaner. This is because it these videos the background is mostly dark with a few bright regions (which implies a highly coherent subspace) and the location of the errors is highly correlated (the location of an object in consecutive frames is very similar). In contrast to other optimization methods [5-12,18,19], we make no assumptions about coherence or the distribution of the sparse errors, and so this does not affect our results.



Figure 2: Sparse (foreground) plus low-rank (background) decomposition of some video frames from the Wallflower [23] and I2R [24] datasets using R2PCA and RPCA-ALM [18,19].