

## ROBUST LOW-RANK MATRIX COMPLETION BY RIEMANNIAN OPTIMIZATION\*

LÉOPOLD CAMBIER<sup>†</sup> AND P.-A. ABSIL<sup>‡</sup>

**Abstract.** Low-rank matrix completion is the problem where one tries to recover a low-rank matrix from noisy observations of a subset of its entries. In this paper, we propose RMC, a new method to deal with the problem of *robust* low-rank matrix completion, i.e., matrix completion where a fraction of the observed entries are corrupted by non-Gaussian noise, typically outliers. The method relies on the idea of smoothing the  $\ell_1$  norm and using Riemannian optimization to deal with the low-rank constraint. We first state the algorithm as the successive minimization of smooth approximations of the  $\ell_1$  norm, and we analyze its convergence by showing the strict decrease of the objective function. We then perform numerical experiments on synthetic data and demonstrate the effectiveness on the proposed method on the Netflix dataset.

**Key words.** low-rank matrix completion, Riemannian optimization, outliers, smoothing techniques,  $\ell_1$  norm, nonsmooth, fixed-rank manifold

**AMS subject classifications.** 65K99, 90C99

**DOI.** 10.1137/15M1025153

**1. Introduction.** The problem of low-rank matrix completion has drawn significant interest in the past decade. It can be used as a building block for recommender systems, where one wants to predict user ratings based on partial ratings using collaborative filtering [4], in reconstructing a three-dimensional path of particles from only partial observation using a fixed camera [18], in sensor network localization [14, 28, 26], or in image inpainting where low-rank completion is used as a way to reconstruct a damaged image [27].

Low-rank matrix completion consists of recovering a rank- $r$  matrix of size  $m \times n$  (with  $r \ll \min(m, n)$ ) from only a fraction (typically  $\mathcal{O}(r(m+n))$  or  $\mathcal{O}(r(m+n) \log(m+n))$ ) of its entries. Denoting by  $\Omega$  the set of observed entries, the problem can be stated as

$$(1) \quad \begin{aligned} \min_{X \in \mathbb{R}^{m \times n}} \quad & \text{rank}(X) \\ \text{subject to} \quad & \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M), \end{aligned}$$

where  $\mathcal{P}_\Omega : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_\Omega^{m \times n} : X \rightarrow \mathcal{P}_\Omega X$  is the orthogonal projector onto the space  $\mathbb{R}_\Omega^{m \times n}$  of  $m \times n$  matrices with zero-entries on  $\bar{\Omega} = \{(i, j) | 1 \leq i \leq m, 1 \leq j \leq n\} \setminus \Omega$ :  $\mathcal{P}_\Omega(X)_{ij} = X_{ij}$  if  $(i, j) \in \Omega$  and 0 otherwise. Finally,  $M$  is the matrix containing the known entries (with values known only on  $\Omega$ ). This problem, however, is now well known to be NP-hard [13].

\*Received by the editors June 15, 2015; accepted for publication (in revised form) March 16, 2016; published electronically October 27, 2016. This paper presents research results of the Belgian Network DYSCO (Dynamical Systems, Control, and Optimization), funded by the Interuniversity Attraction Poles Programme initiated by the Belgian Science Policy Office.

<http://www.siam.org/journals/sisc/38-5/M102515.html>

<sup>†</sup>Institute for Computational & Mathematical Engineering, Stanford University, Stanford, CA 94305 (lcambier@stanford.edu).

<sup>‡</sup>ICTEAM Institute, Université catholique de Louvain, B-1348 Louvain-la-Neuve, Belgium (absil@inma.ucl.ac.be). This author’s work was supported by “Communauté française de Belgique - Actions de Recherche Concertées” and by FNRS under grant PDR T.0173.13.

Candès and Recht [10] stated the problem as

$$\begin{aligned} & \min_{X \in \mathbb{R}^{m \times n}} \|X\|_* \\ & \text{subject to } \mathcal{P}_\Omega(X) = \mathcal{P}_\Omega(M), \end{aligned}$$

where  $\|\cdot\|_*$  is the nuclear norm  $\|X\|_* = \sum_{k=1}^{\min(m,n)} \sigma_k(X)$  with  $\sigma_k(X)$  the  $k$ th singular value of  $X$ . The authors proved that in the context of *exact* low-rank matrix completion, this formulation recovers the original underlying matrix under some mild assumptions.

Another way to approach low-rank matrix completion is the following. Assume the rank  $r$  of the target matrix is known in advance (which does make sense in a lot of applications such as computer vision, where the rank is often related to the dimension of the space). In this case, the problem can be stated as

$$(2) \quad \min_{X \in \mathcal{M}_r} \|\mathcal{P}_\Omega(X - M)\|_{\ell_2},$$

where

$$\mathcal{M}_r = \{X \in \mathbb{R}^{m \times n} : \text{rank}(X) = r\}$$

and where  $\|\cdot\|_{\ell_2} = \|\cdot\|_F$  is the  $\ell_2$  or Frobenius norm.<sup>1</sup> Intuitively, this problem seeks the matrix  $X$  of rank  $r$  that best fits the given data. The main advantage is that it is robust to Gaussian additive noise, in the sense that a small Gaussian additive noise still allows recovery of the underlying low-rank matrix with an error proportional to the noise level [19]. This problem and other similar formulations have been addressed by different authors.

Vandereycken [29] takes advantage of the fact that  $\mathcal{M}_r$  is a smooth Riemannian manifold to apply recent optimization algorithms [1] to efficiently solve the problem. Our method will use the exact same tools.

All formulations that rely on the Frobenius norm as in (2) suffer from one drawback: even though they are robust to additive Gaussian noise, they are not well suited to recover the underlying low-rank matrix when the noise becomes sufficiently far from Gaussian. Here we focus on the situation where only a few of the observed entries, termed *outliers*, are perturbed; that is,

$$(3) \quad M = M_0 + S,$$

where  $M_0$  is the unperturbed data matrix of rank  $r$  and  $S$  is a sparse matrix. For instance, consider recovering the best rank-1 approximation of the following matrix:

$$M_x = \begin{pmatrix} 2 & -1+x \\ 4 & -2 \end{pmatrix}.$$

If  $x = 0$ , this matrix is rank 1 since, for instance,

$$M_0 = \begin{pmatrix} 2 & -1 \\ 4 & -2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \cdot \begin{pmatrix} 2 & -1 \end{pmatrix}.$$

But when  $x \neq 0$ , this is not the case, and finding the rank-1 matrix that minimizes the  $\ell_2$  error leads to fundamentally different solutions.

---

<sup>1</sup>In this paper, we use the notation  $\|\cdot\|_{\ell_2}$  for the Frobenius norm to emphasize the difference from the  $\ell_1$  norm,  $\|\cdot\|_{\ell_1}$ .

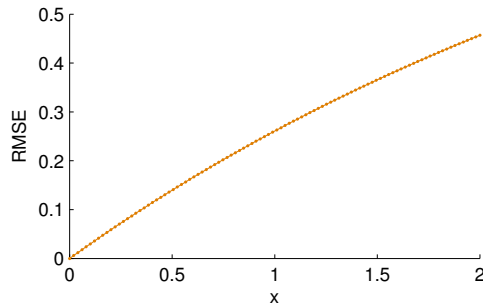


FIG. 1. The error between the recovered matrix and the original one when minimizing the  $\ell_2$  norm, as a function of  $x$ . We can clearly observe that when using an  $\ell_2$  loss function, even a small perturbation on a single entry of the matrix can lead to a large error with respect to the original matrix.

To observe this, we can simply compute, for each  $x$ , the rank-1 SVD of  $M_x$  (which is the solution that an  $\ell_2$  method minimizing  $\|M_x - X\|_{\ell_2}$  would return) and then compute the root mean square error (RMSE) with respect to the original matrix  $M_0$ . This is depicted in Figure 1, and we can see that the error starts to grow as soon as  $x \neq 0$ .

However, if the method was able to identify that only the top right entry of  $M_x$  is corrupted by noise, then it would be able to recover  $M_0$  exactly by removing the top right entry from the mask  $\Omega$  and performing low-rank matrix completion. More generally, in the problem of completing from its known entries  $\Omega$  a matrix  $M$  generated as in (3), the ability to detect the outliers in  $\mathcal{P}_\Omega(M)$  (i.e., the entries affected by the sparse matrix  $\mathcal{P}_\Omega(S)$ ) and remove those entries from the mask  $\Omega$  would open the way for an exact recovery of the rank- $r$  matrix  $M_0$ .

**1.1. Previous work.** Matrix completion in the presence of outliers has been considered in several papers.

Chen et al. [12] studied the problem of low-rank matrix completion where a large number of columns are arbitrarily corrupted. They showed that only a small fraction of the entries are needed in order to recover the low-rank matrix with high probability, without any assumptions on the location or the amplitude of the corrupted entries.

Both Li [22] and Chen et al. [11] studied a harder problem, when a constant fraction of the entries (not the columns) of the matrix are outliers. They studied what conditions need to be imposed in order for the convex optimization problem

$$\begin{aligned} \min \quad & \gamma \|X\|_* + \|E\|_{\ell_1} \\ \text{subject to} \quad & \mathcal{P}_\Omega(X + E) = \mathcal{P}_\Omega(M) \end{aligned}$$

to exactly recover the underlying low-rank matrix (with  $\|\cdot\|_*$  the nuclear norm). Basically, they showed that there exist universal constants such that with overwhelming probability the solution of the problem is equal to  $M$  on the mask  $\Omega$ . In the close context of low-rank principal component analysis, Candès et al. [8] was also able to solve the same problem. The advantage of such an algorithm is that it is convex and can then be analyzed thoroughly.

This robust formulation has been improved to deal with Gaussian noise [15], leading to the following convex optimization problem (convex robust matrix completion,

or CRMC):

$$(4) \quad \begin{aligned} \min \quad & \lambda \|X\|_* + \gamma \|E_1\|_{\ell_1} + \frac{1}{2} \|E_2\|_{\ell_2}^2 \\ \text{subject to} \quad & \mathcal{P}_\Omega(X + E_1 + E_2) = \mathcal{P}_\Omega(M). \end{aligned}$$

He and colleagues [16, 17] developed a robust version of the GROUSE algorithm [3], named GRASTA, which aims at solving the problem of *robust* subspace tracking. Their algorithm can be cast to solve problems formulated as

$$\begin{aligned} \min \quad & \|\mathcal{P}_\Omega(S)\|_{\ell_1} \\ \text{subject to} \quad & \mathcal{P}_\Omega(UV + S) = \mathcal{P}_\Omega(M), \\ & U \in \text{Gr}(m, r), \\ & V \in \mathbb{R}^{r \times n}, \end{aligned}$$

where  $\text{Gr}(m, r)$  is the Grassman manifold, i.e., the set of linear  $r$ -dimensional subspaces of  $\mathbb{R}^m$ . GRASTA tackles this problem by first building the augmented Lagrangian problem, and then solves it by alternating between  $V$ , its dual variables, and  $U$  and by performing steepest descent on the Grassman manifold. The advantage of their algorithm is that it is designed to tackle the problem of *online* subspace estimation from incomplete data; hence it can also be casted to solve online low-rank matrix completion where we observe one column of the matrix  $M$  at a time.

Nie and colleagues [24, 25] solved a slightly more general problem where all norms become arbitrary  $p$ -norms:

$$\min \lambda \|X\|_{S_p}^p + \|\mathcal{P}_\Omega(X - M)\|_{\ell_p}^p,$$

where  $\|X\|_{S_p}^p = \sum_{i=1}^{\min(m,n)} \sigma_i^p(X)$  and  $\|X\|_{\ell_p}^p = \sum_{i=1, j=1}^{m,n} |X_{ij}|^p$ . The algorithm used to solve this nonconvex program (when  $p < 1$ ) is, again, an augmented Lagrangian method. We were unfortunately unable to obtain or write an efficient implementation of this algorithm since it requires the storage of the full  $m \times n$  matrix, as well as SVD of full matrices of this size. This formulation, however, is efficient for moderate size problems.

Yan, Yang, and Osher [31] solved  $\ell_2$  problems of the form

$$\min_{X \in \mathcal{M}_r} \|\mathcal{P}_\Omega(X - M)\|_{\ell_2},$$

where the mask  $\Omega$  is adapted at each iteration to remove the suspected outliers. The idea is to first solve the problem with the original mask  $\Omega$ , detect outliers, adapt the mask, and then solve the problem again until convergence. Intermediate problems are handled using RTRMC [5].

Yang, Fend, and Suykens [32] studied the problem of robust low-rank matrix completion using a nonconvex loss-function. They solve the following problem:

$$\min_{X \in \mathbb{R}^{m \times n}: \text{rank}(X) \leq r} \frac{\sigma^2}{2} \sum_{(i,j) \in \Omega} \left( 1 - \exp\left(-\frac{(X_{ij} - M_{ij})^2}{\sigma^2}\right) \right),$$

where the rankconstraint is relaxed using the now standard nuclear-norm heuristic.

Finally, Klopp, Lounici, and Tsybakov [20] studied the optimal reconstruction error in the case of matrix completion, where the observations are noisy and column-wise or elementwise corrupted and where the only piece of information needed is a bound on the matrix entries. They provided a range of (optimal) estimators to solve such problems with guarantees.

**1.2. Contribution.** In this paper, we consider low-rank matrix completion in the presence of outliers, assuming as in, e.g., [31] that the rank  $r$  is known in advance, which naturally leads to the formulation

$$\min_{X \in \mathcal{M}_r} \|\mathcal{P}_\Omega(X - M)\|_{\ell_p}.$$

It remains to choose  $p$ . The choice  $p = 0$ , which results in maximizing the number of exactly recovered entries over the mask  $\Omega$ , seems natural when only a sparse noise is present, but this discontinuous objective function is unwieldy, and, moreover, it is inadequate in the presence of an additional dense (i.e., nonsparse) noise. The choice  $p = 2$ , as in [29, 6], would be adequate for Gaussian additive noise, but its mean square nature makes it excessively sensitive to the outliers. We opt for the middle ground, namely  $p = 1$ . Its well-known sparsity-inducing property lets us expect exact recovery when the noise consists of just a few outliers. We will see in the numerical experiments that this is indeed the case.

The choice  $p = 1$  leaves us with a nonsmooth objective function. We handle this difficulty by replacing successively the  $\ell_1$  norm by increasingly accurate smooth approximations thereof. As in [31], the resulting minimization problems over the fixed-rank manifold  $\mathcal{M}_r$  are tackled by Riemannian optimization techniques. However, while the  $\ell_2$  formulation in [31] enables a variable projection strategy that yields an optimization problem on the Grassmann manifold [6], our smoothed  $\ell_1$  objective functions do not lend themselves to this approach. Following the way paved in [29], we resort instead to a conjugate gradient (CG) scheme on the fixed-rank manifold  $\mathcal{M}_r$  viewed as an embedded Riemannian submanifold of  $\mathbb{R}^{m \times n}$ . Numerical experiments confirm that the resulting algorithm is particularly efficient in the case where a few percent of the entries are largely corrupted by non-Gaussian noise.

Compared to the methods described by [8] or [15], for instance, our method requires at least an estimate of the target rank. In some applications, as in computer vision, the target rank can be known in advance. It can also often be estimated, and then adjusted according to the result. In comparison with the nuclear-norm formulation that requires full SVD factorizations, iterating on low-rank manifolds yields lower time and space complexities.

We also point out that our problem formulation, like that of GRASTA, involves an  $\ell_1$  objective function and a fixed-rank constraint. However the algorithms are fundamentally different: GRASTA proceeds one column at a time and combines gradient descent on the Grassmannian and the alternating direction method of multipliers, while our proposed method applies to the whole objective function a CG scheme on a fixed-rank manifold.

Our algorithm can also efficiently handle regularization, and in practice we solve

$$\min_{X \in \mathcal{M}_r} \|\mathcal{P}_\Omega(X - M)\|_{\ell_1} + \lambda \|\mathcal{P}_{\bar{\Omega}}(X)\|_{\ell_2}^2.$$

The regularization factor appears to be quite useful in applications to limit overfitting. This is justified by the work of Boumal and Absil [6], where regularization is crucial to handling real datasets. The algorithm also scales well with the size of the problem and stays very efficient when the amplitude of the outliers increases. We finally implemented a simple MATLAB version of the algorithm, able to solve problems of size  $50000 \times 50000$  of rank 10 in a matter of minutes on a regular desktop computer.

This paper is divided in the following way. In section 2 we remind the reader of the essential tools of optimization on manifolds. We introduce our algorithm in

section 3. We study its convergence in section 4, while we perform numerical experiments, both synthetic ones and real-life applications, in section 5. Conclusions are drawn in section 6.

**2. Optimization on manifolds and the low-rank matrix manifold.** The *low-rank matrix manifold*

$$\mathcal{M}_r = \{X \in \mathbb{R}^{m \times n} : \text{rank}(X) = r\},$$

where  $r \leq \min(m, n)$ , is known to be a smooth manifold embedded in  $\mathbb{R}^{m \times n}$  of dimension  $r(m + n - r)$  [21, 30]. Hence, optimization techniques presented in [1] can be applied to solve smooth optimization problems where constraints are formulated using  $\mathcal{M}_r$ .

There are several ways of describing a matrix  $X \in \mathcal{M}_r$  [2]. In this paper, we will use the very natural SVD-like representation

$$(5) \quad X = U \Sigma V^\top,$$

where  $U \in \mathbb{R}^{m \times r}$  and  $V \in \mathbb{R}^{n \times r}$  are matrices with orthogonal columns (i.e.,  $U^\top U = I_r$  and  $V^\top V = I_r$ ) and  $\Sigma \in \mathbb{R}^{r \times r}$  is a diagonal full-rank matrix. This formulation requires  $r(m + n + 1) \approx r(m + n - r)$  storage capacity and has the advantage of having two orthogonal matrices.

Each vector  $\dot{X}$  belonging to the tangent space  $T_X \mathcal{M}_r$  of  $\mathcal{M}_r$  at  $X$  has a unique representation  $(\dot{U}, \dot{\Sigma}, \dot{V})$  such that [29]

$$(6) \quad \begin{aligned} \dot{X} &= U \dot{\Sigma} V^\top + \dot{U} V^\top + U \dot{V}^\top, \\ U^\top \dot{U} &= 0 \text{ and } V^\top \dot{V} = 0. \end{aligned}$$

This formulation also requires  $r(m + n + 1) \approx r(m + n - r)$  storage capacity.

Given a vector  $Z$  in the ambient space  $\mathbb{R}^{m \times n}$ , its projection on the tangent space  $T_X \mathcal{M}_r$  can be computed [29] and is given by  $P_{T_X \mathcal{M}_r}(Z)$ , defined as

$$P_{T_X \mathcal{M}_r} : \mathbb{R}^{m \times n} \rightarrow T_X \mathcal{M}_r : Z \rightarrow P_U Z P_V + P_U^\perp Z P_V + P_U Z P_V^\perp,$$

with  $P_U = U U^\top$  and  $P_U^\perp = I - U U^\top$ ,  $P_V$  and  $P_V^\perp$  being defined in the same way.

Using the tangent space representation, a basic identification yields the following representation for the (orthogonal) projection of an ambient vector  $Z$  onto  $T_X \mathcal{M}_r$ :

$$(7) \quad \dot{\Sigma} = U^\top Z V, \quad \dot{U} = (I - U U^\top) Z V, \quad \dot{V} = (I - V V^\top) Z^\top U.$$

The algorithm we will describe requires being able to move along directions on the manifold. It is now well established (e.g., in [1]) that this can be cheaply achieved using a *retraction* instead of the expensive exponential map, for instance, while keeping all convergence guarantees. We decided to use the projective retraction [2]: given a vector  $\dot{X} \in T_X \mathcal{M}_r$ , it finds  $Y \in \mathcal{M}_r$  such that

$$Y = R_X(\dot{X}) = \underset{Y \in \mathcal{M}_r}{\text{argmin}} \|X + \dot{X} - Y\|_F.$$

The solution of this minimization problem is known to be the rank- $r$  SVD of  $X + \dot{X}$  (Eckart–Young theorem). Assuming  $X$  is given as (5) and  $\dot{X}$  as (6), it is possible to compute it efficiently [29, 2].

Because  $\mathcal{M}_r$  is embedded in  $\mathbb{R}^{m \times n}$ , a suitable vector transport (i.e., a mapping from the tangent space at some point to the tangent space at another point) is simply the projection of the ambient version of the original vector in the tangent space at the new point [29].

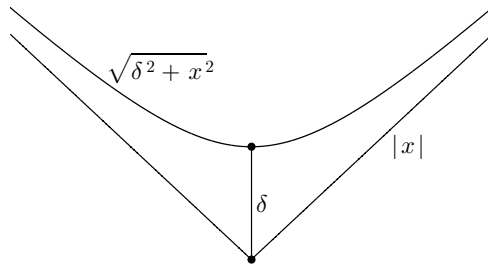


FIG. 2. Illustration of both  $|x|$  and  $\sqrt{\delta^2 + x^2}$ .

### 3. Riemannian optimization and smoothing techniques.

**3.1. The main idea.** As explained earlier, the problem we aim to solve is

$$(8) \quad \min_{X \in \mathcal{M}_r} \|\mathcal{P}_\Omega(X - M)\|_{\ell_1} + \lambda \|\mathcal{P}_{\bar{\Omega}}(X)\|_{\ell_2}^2,$$

with the following interpretation: we ought to find a low-rank matrix  $X$  that fits the data  $M$  in the  $\ell_1$  sense on the mask  $\Omega$ . On the remaining entries in  $\bar{\Omega}$ , we have a small confidence  $\lambda$  (typically between 0 and  $10^{-5}$ – $10^{-3}$ , even though this is application-dependent) that the value should be zero; hence we minimize the  $\ell_2$  error between  $X$  and 0 on  $\bar{\Omega}$ .<sup>2</sup> This is motivated by previous studies of [6], where regularization was especially useful in dealing with real datasets. Note that the reason for the use of the  $\ell_2$  norm in the regularization term is twofold: first, the  $\ell_2$  norm will allow significant simplifications to be detailed in the forthcoming sections. Second, we can observe outliers on  $\Omega$ , so the  $\ell_1$  norm makes sense there; but we obviously cannot observe outliers on  $\bar{\Omega}$ , so it does not seem necessary to use an  $\ell_1$  norm there, and an  $\ell_2$  norm should be more suitable.

The obvious main drawback of using (8) is that it is nondifferentiable. To remedy this problem, we decided to use smoothing techniques in order to make the objective differentiable. The idea is that, for a small  $\delta > 0$ , the function

$$\sum_{(i,j) \in \Omega} \sqrt{\delta^2 + (X_{ij} - M_{ij})^2}$$

is a smooth approximation of

$$\|\mathcal{P}_\Omega(X - M)\|_{\ell_1},$$

as depicted in Figure 2. The idea is thus to solve the optimization problem

$$(9) \quad \min_{X \in \mathcal{M}_r} \sum_{(i,j) \in \Omega} \sqrt{\delta^2 + (X_{ij} - M_{ij})^2} + \lambda \sum_{(i,j) \in \bar{\Omega}} X_{ij}^2$$

for decreasing values of  $\delta$ .

To solve a problem in the form

$$\min_{x \in \mathcal{M}} f(x),$$

<sup>2</sup>Note that this assumes that the entries in the matrix have a mean equal to zero.

where  $\mathcal{M}$  is a smooth Riemannian manifold and where  $f$  is smooth, there now exist many different techniques such as Riemannian CG or trust-region algorithms [1]. We have opted for the CG approach, as it appears to be more precise and efficient when  $\delta$  becomes small.

**3.2. The objective function and its gradient.** Using a first-order algorithm like CG requires the computation of the cost function and the (Riemannian) gradient.

Taking a look at (9), one may think that just evaluating the cost would require  $\mathcal{O}(mn)$  operations, since we need to evaluate  $X$  over *both*  $\Omega$  and  $\bar{\Omega}$ . Actually, as pointed out in [5], this is not the case, since

$$\|\mathcal{P}_{\bar{\Omega}}(X)\|_{\ell_2}^2 = \|X\|_{\ell_2}^2 - \|\mathcal{P}_{\Omega}(X)\|_{\ell_2}^2,$$

and, because we store  $X$  using the factorization  $X = U\Sigma V^\top$ , we can compute  $\|X\|_{\ell_2}^2$  easily thanks to the invariance of the Frobenius norm to orthogonal changes of basis:

$$\|X\|_{\ell_2}^2 = \|\Sigma\|_{\ell_2}^2,$$

which requires  $\mathcal{O}(r)$  operations. We can then rewrite the cost function of (9) as

$$(10) \quad f_\delta(X) = \sum_{(i,j) \in \Omega} \left( \sqrt{\delta^2 + (X_{ij} - M_{ij})^2} - \lambda X_{ij}^2 \right) + \lambda \|X\|_{\ell_2}^2.$$

Hence, computing the cost function requires  $\mathcal{O}(|\Omega| + r)$  operations, after having evaluated the product  $U\Sigma V^\top$  on the mask  $\Omega$ .

We can clearly see here the advantage of having added an  $\ell_2$  regularization term: the fact that  $\|X\|_{\ell_2}^2 = \|\Sigma\|_{\ell_2}^2$  is crucial to avoid an  $\mathcal{O}(mn)$  complexity in the computation of the objective function.

To compute the gradient, because  $\mathcal{M}_r$  is embedded in  $\mathbb{R}^{m \times n}$ , we first need to compute the Euclidean gradient of  $f$  at  $X$  and then project it onto  $T_X \mathcal{M}_r$ . The Euclidean gradient is

$$\nabla f_\delta(X) = S + 2\lambda X,$$

where  $S$  is a *sparse matrix* defined as

$$S_{ij} = \begin{cases} \frac{X_{ij} - M_{ij}}{\sqrt{\delta^2 + (X_{ij} - M_{ij})^2}} - 2\lambda X_{ij} & \text{if } (i, j) \in \Omega, \\ 0 & \text{otherwise.} \end{cases}$$

The gradient is thus the sum of a sparse ( $S$ ) and a low-rank ( $2\lambda X$ ) component. Then, projecting it onto  $T_X \mathcal{M}_r$  can be done efficiently thanks to (7) and to the factorization  $X = U\Sigma V^\top$ . Indeed, we have

$$\begin{aligned} \dot{\Sigma} &= U^\top (S + 2\lambda X) V \\ &= U^\top S V + 2\lambda \Sigma, \\ \dot{U} &= (I - U U^\top) (S + 2\lambda X) V \\ &= S V + 2\lambda U \Sigma - U U^\top S V - 2\lambda U \Sigma \\ &= S V - U U^\top S V, \\ \dot{V} &= (I - V V^\top) (S + 2\lambda X)^\top U \\ &= S^\top U + 2\lambda V \Sigma - V V^\top S^\top U - 2\lambda V \Sigma \\ &= S^\top U - V V^\top S^\top U. \end{aligned}$$



Given the fact that  $S$  is sparse, these three terms can be computed efficiently. Note that the addition of the regularization parameter is cheap, since it only requires us to modify the  $S$  matrix and to add a small  $r \times r$  matrix to  $\tilde{\Sigma}$ .

**3.3. The algorithm.** The full algorithm is stated in Algorithm 1; the name RMC stands for “robust matrix completion.” Note that in the following, by outer and inner iteration we mean the  $\delta$  and the CG loop, respectively.

We use a CG algorithm with a Hestenes–Stiefel modified rule (even though, after several experiments, we found that this choice does not really impact the algorithm) and an Armijo backtracking linesearch.

The starting point of the algorithm,  $X^{(0)}$ , can be chosen simply using the rank- $r$  SVD of  $\mathcal{P}_\Omega(M)$ . Suitable values for  $\delta^{(0)}$  (the initial value for the smoothing parameter  $\delta$ ) are application-dependent, but for data  $M$  with values around unity, we use  $\delta^{(0)} = 1$ . Note that this value can have a significant impact on the quality of the final solution. The smoothing parameter is then updated using a geometric rule  $\delta^{(k+1)} = \theta \cdot \delta^{(k)}$ . A quite “aggressive” value of  $\theta = 0.05$  gives good results in our synthetic experiments. In real applications this parameter has to be tuned to find a suitable value. For all experiments,  $\epsilon$  is set to  $10^{-8}$  (but again, this is application-dependent). The stopping criterion of the CG algorithm is set to a maximum of 40 iterations or a gradient norm of  $10^{-8}$ , whichever is reached first.

---

**Algorithm 1** RMC
 

---

**procedure** RMC( $X^{(0)}, \delta^{(0)}, \theta, \epsilon, \lambda$ )

$f^{(0)} \leftarrow \infty$

$k \leftarrow 0$

$\delta^{(1)} \leftarrow \delta^{(0)}$

$e \leftarrow \infty$

**while**  $e \geq \epsilon$  **do**

Solve

$$(11) \quad X^{(k+1)} = \operatorname{argmin}_{X \in \mathcal{M}_r} \sum_{(i,j) \in \Omega} \sqrt{(\delta^{(k+1)})^2 + (X_{ij} - M_{ij})^2} + \lambda \|\mathcal{P}_\Omega(X)\|_{\ell_2}^2$$

using Riemannian CG algorithm starting from  $X^{(k)}$ .

$f^{(k+1)} \leftarrow f_{\delta^{(k+1)}}(X^{(k+1)})$

$e \leftarrow f^{(k)} - f^{(k+1)}$

$k \leftarrow k + 1$

$\delta^{(k+1)} \leftarrow \delta^{(k)} \cdot \theta$

**end while**

**end procedure**

---

**4. Convergence analysis.** This section provides a basic convergence analysis of the RMC algorithm. Its goal is mostly to justify the stopping criterion of the outer loop, i.e., that the algorithm will terminate at some point (assuming exact arithmetic at least). Let  $f_\delta$  be defined as in (10).

**THEOREM 1** (strict decrease). *If the sequence of iterates  $\{X^{(0)}, X^{(1)}, X^{(2)}, \dots\}$  is produced by Algorithm 1 with  $\theta < 1$ , defining  $f^{(k)} = f_{\delta^{(k)}}(X^{(k)})$ , we have*

$$f^{(0)} > f^{(1)} > \dots > f^{(k)} > \dots$$

*Proof.* At iteration  $k$ , the CG algorithm returns a feasible solution  $X^{(k)} \in \mathcal{M}_r$ , associated with  $\delta^{(k)}$ . It is easy to see that  $X^{(k)}$  stays a feasible point for the next step (since it belongs to  $\mathcal{M}_r$ ), and that

$$f_{\delta^{(k)}}(X^{(k)}) > f_{\delta^{(k+1)}}(X^{(k)}).$$

This follows from  $\delta^{(k+1)} = \theta \cdot \delta^{(k)} < \delta^{(k)}$  and the expression (10) of  $f_\delta(X)$ . Then, because CG is a descent direction,

$$f_{\delta^{(k+1)}}(X^{(k)}) \geq f_{\delta^{(k+1)}}(X^{(k+1)}).$$

The claim follows from these two inequalities. □

Now, it is also easy to notice that, for all  $\delta \geq 0$  and for all  $X \in \mathcal{M}_r$ ,

$$f_\delta(X) \geq f(X).$$

Hence, defining

$$f^* = \inf_{X \in \mathcal{M}_r} f(X) \geq 0,$$

we observe that the sequence of iterates  $\{f_k\}_{k=1}^K$  is monotonically decreasing and bounded below by  $f^*$ .

This conclusion justifies the stopping criterion of the algorithm saying that it stops after iteration  $k$  if the difference between  $f^{(k)}$  and  $f^{(k+1)}$  is below some threshold  $\epsilon$ : the algorithm will terminate at some point. We emphasize the fact that this does not prove that the algorithm converges toward a global minimum.

**5. Numerical experiments.** In this section, we first apply RMC on synthetic problems. We then perform experiments on the Netflix dataset to show how the algorithm behaves in this situation.

**5.1. Synthetic problems.** In all experiments, synthetic data are created in the following way: we build  $U \in \mathbb{R}^{m \times r}$  and  $V \in \mathbb{R}^{r \times n}$  with independent and identically distributed (i.i.d.) Gaussian entries such that their product  $M = UV$  is filled with zero-mean and unit-variance nonindependent Gaussian entries. We then sample  $k = \rho r(m + n - r)$  entries uniformly at random, where  $\rho$  is the oversampling factor.

In some cases, we will add some non-Gaussian noise on part the observed entries to create outliers. To do so, we add one realization of the random variable

$$\mathcal{O} = \mathcal{S}_{\pm 1} \cdot \mathcal{N}(\mu, \sigma^2),$$

where  $\mathcal{S}_{\pm 1}$  is a random variable with equal probability of being equal to  $+1$  or  $-1$ , while  $\mathcal{N}(\mu, \sigma^2)$  is a Gaussian random variable of mean  $\mu$  and variance  $\sigma^2$ . Outliers are always created uniformly at random with a probability of 5%.

Different factors affect the task of matrix completion. Obviously, the size of the problem matters, and we will try to tackle large enough problems to show that our algorithm scales well. The oversampling factor  $\rho$  should also be greater than 1, and in the following experiments it will be fixed to either 4 or 5.

Let us denote by  $M_0$  the original low-rank matrix, without any outliers. We will monitor how the RMSE, defined as the error on *all* the entries between  $X$  and the *original* matrix  $M_0$ , i.e.,

$$RMSE(X, M_0) = \sqrt{\frac{\sum_{i=1, j=1}^{m, n} (X_{ij} - M_{0,ij})^2}{mn}},$$

decreases. Since we have access to the factorization of both  $X$  and  $M_0$ , this can be computed efficiently [6]. A decrease toward zero is what we expect from an RMC method, since our goal is to recover exactly (up to numerical errors) the original low-rank matrix, even in the presence of outliers.

We decided to compare RMC (Algorithm 1) to AOPMC [31], GRASTA [16, 17], CRMC [15], and the  $\ell_2$  method RTRMC [6].

For RMC, the maximum number of CG iterations (the inner loop) is set to 40 with a gradient tolerance of  $10^{-8}$ . We use  $\delta_0 = 1$ , as well as  $\theta = 0.05$ .

For AOPMC, we use the default settings with a maximum of 20 trust-region iterations at each outer iteration (i.e., when the mask  $\Omega$  is fixed, with potentially some outliers removed) but a maximum of 20 iterations for the tCG algorithm (see [5] for further information). Note that we use the code of the authors,<sup>3</sup> but because we know the number of outliers, we decided to provide it to AOPMC. Otherwise, we would need to run the algorithm several times to guess the number of outliers. Also note that AOPMC automatically adjusts the number of iterations if it seems that the convergence is too bad. We did not change this option, so this may explain why the algorithm sometimes does more than 20 iterations between each update of the mask  $\Omega$ .

Regarding GRASTA, we also use the implementation provided by the authors,<sup>4</sup> with the default settings, but we had some trouble running the algorithm on large  $50000 \times 50000$  problems, since it was not even able to perform two complete sweeps over the data in less than 10 minutes (while RMC terminates in about 5 minutes in this situation). This is due to the fact that the algorithm operates one column at a time, and since each rank-1 update of the  $U$  matrix takes about 0.01–0.02 second (using a standard but efficient MATLAB implementation), one loop over the 50000 columns takes more than 10 minutes. Note that this can be easily understood, since the original goal of GRASTA is not to perform “batch” matrix completion, but rather *online* subspace tracking. For the  $5000 \times 5000$  case, everything goes well and the algorithm converges in a decent amount of time. The only modification made to the algorithm was to change the initial point that was set, like the other algorithms, to the left matrix of the rank- $r$  SVD of the matrix  $\mathcal{P}_\Omega(M)$  (instead of some complete random initialization). Note that this does not have a significant influence on the convergence of the algorithm.

For CRMC, we obtained the code directly from the authors. It uses the PROPACK toolbox to compute the SVD. All the computational burden comes from the required *dense* SVD (required to compute the gradient step) that has to be performed at every step. As we will show in the next experiment, this significantly slows down the algorithm. CRMC has two essential parameters (see equation (4)), namely  $\gamma$  and  $\lambda$ . In our experiments, we noticed a sharp increase in the solution quality for  $\lambda = 10^{-6}$  and  $\gamma = 4.5 \cdot 10^{-3}$ . Note that since RMC knows the rank of the target low-rank matrix, we decided to provide it to CRMC: it significantly speeds up the convergence by only computing thin SVDs.

Finally, we compared these four  $\ell_1$  methods with an efficient  $\ell_2$  method. This aims at showing that  $\ell_2$  methods do not perform well in the presence of outliers, even for very moderate values. We decided to use RTRMC [6] version 3.1, with all the default parameters, including the regularization parameter  $\lambda$ , set to 0. Note that higher values do not lead to better solutions in the presence of outliers.

<sup>3</sup>Available at <https://binary-matching-pursuit.googlecode.com/files/AOPMCv1.zip>.

<sup>4</sup>Available at <https://sites.google.com/site/hejunzz/grasta>, version 1.2.0.

In Figures 3–6, the large dots indicate a change in the outer iteration: in RMC it indicates a decrease in  $\delta$ , while it indicates an update of the mask  $\Omega$  in AOPMC. The dash-dot (with small dots) line is used to distinguish RTRMC from other methods.

*Remark 1.* We point out that we tried, without success, to speed up the RMC and AOPMC algorithms by terminating the inner loop sooner, when the decrease from iteration to iteration was small enough. It appears that solving the inner problems with a good enough quality (i.e., a small enough gradient norm) is crucial for the convergence of both algorithms toward the exact underlying low-rank matrix. For this reason, we left the gradient tolerance set to  $10^{-8}$ . This explains the flat regions at the end of each iteration of the outer loop on the following figures.

*Remark 2.* Note that the experiments were run on quite large  $50000 \times 50000$  matrices in order to show that the algorithm presented here scales well on large matrices. Similar experiments were also run on  $500000 \times 500000$  matrices of rank 10: the behavior of RMC was very stable, and the running time appeared to be linear in the size of the matrix, i.e., in  $\mathcal{O}(m+n)$ . They are omitted for simplicity. Results on smaller matrices are qualitatively the same, except for the time the algorithm takes to reach the optimal value.

All experiments were run on a 6-core Intel Xeon CPU E5-1650 v2 at 3.50GHz with 64 GB of RAM using MATLAB R2014a on a 64-bit Linux machine. We used the Manopt toolbox [7] (version 1.0.7) to handle the optimization part of the RMC algorithm with the `fixedrankembeddedfactory` manifold factory and the default CG method.

**5.1.1. Perfect low-rank matrix completion.** As a sanity check, we test all the methods on the very simple *perfect matrix completion* (matrix completion without any noise or outliers) problem using a  $5000 \times 5000$  matrix of rank 10. Results are depicted in Figure 3. All methods eventually successfully recover the original matrix: the RMSE is driven toward zero, as expected.

GRASTA and CRMC are observed to be much slower than the other three methods. In both methods, the first iterate is computed after 2–5 seconds (explaining the fact that the curves start only later). For GRASTA, as explained earlier, this is due to the fact that it operates one column at a time. In this particular case, it takes GRASTA approximately 40 seconds to reach an RMSE of  $10^{-8}$ . For CRMC, the dense SVDs make the algorithm so slow that it cannot even decrease the RMSE by a factor 10 in 5 minutes.

**5.1.2. Low-rank matrix completion with outliers.** Given a  $500 \times 500$  matrix for which we observe the entries uniformly at random with an oversampling  $\rho$  of 5, we perturbed 5% of the observed entries by adding to them some non-Gaussian noise to create outliers. This problem would be cumbersome to solve with an  $\ell_2$  method because of the high weights the outliers would have in the objective function (as depicted by the result of RTRMC in the following plots).

When running our algorithms, we obtain the results depicted in Figure 4(a) for outliers created using  $\mu = \sigma = 0.1$  and in Figure 4(b) using  $\mu = \sigma = 1$ .

We can see that all the  $\ell_1$  methods manage to successfully solve this problem. CRMC is quite slow compared to the other methods due to the expensive dense SVD required at each iteration. Still, it is able to reduce the error by three orders of magnitude. This is in contrast with the  $\ell_2$  method RTRMC. We observe that RTRMC follows exactly the same convergence as the first outer iteration of AOPMC, but it then stops. The high weight of the outliers in the  $\ell_2$  objective function prevent any

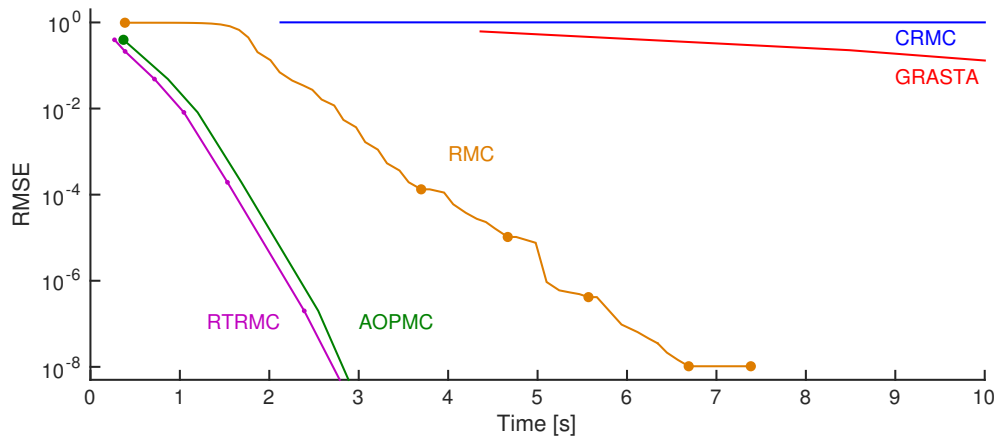


FIG. 3. *Perfect low-rank matrix completion: Low-rank matrix completion of a rank-10  $5000 \times 5000$  matrix observed with an oversampling of 5. The decrease in the objective function from iteration to iteration is clear. In this example, RMC struggles to significantly decrease the RMSE at the end since the function becomes less and less differentiable near the “kinks” of the absolute values, where the solution is located.*

improvement in the solution quality. We can see that, in this case, the strength of the outliers does not have a significant impact: it is moderate enough so that all the  $\ell_1$  methods successfully solve the problem.

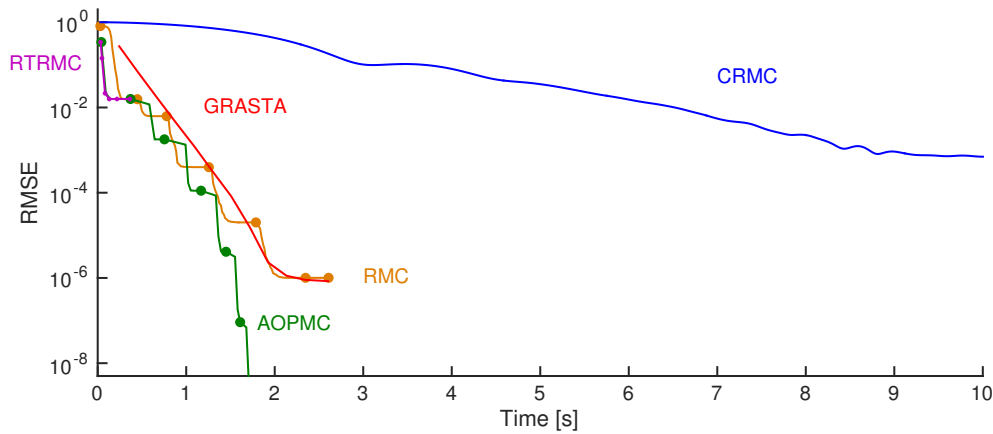
We then run the same experiment on larger  $50000 \times 50000$  matrices, with still 5% of outliers. Figures 5(a) and 5(b) illustrate the results of these experiments, with  $\mu = \sigma = 1$  and  $\mu = \sigma = 5$ , respectively, using an oversampling of 5.

We can see that both AOPMC and RMC solve the first problem well. Both GRASTA and CRMC, on the other hand, do not converge in a decent amount of time. We also observe that RMC stays very robust when the strength of the outliers increases, while AOPMC starts to have important difficulties in the second experiment. The robustness of RMC is most likely due to the asymptotic linear behavior of the cost function, even in the first iterations: AOPMC, on the contrary, needs to first solve an  $\ell_2$  problem. If this problem is too hard, the first outer iteration will lead to such a bad solution that the algorithm will not eventually converge. Finally, note that RTRMC, the  $\ell_2$  method, simply does not converge in the second case.

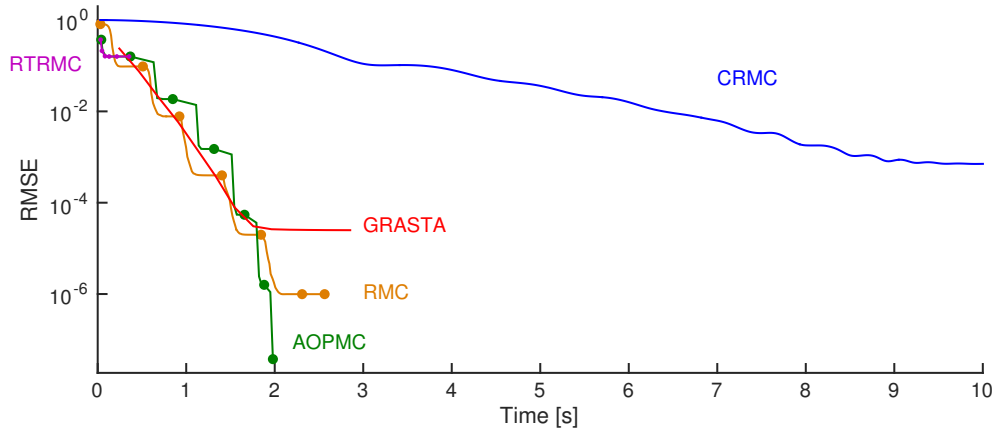
Note that this is a quite extreme experiment, in the sense that the outliers have a mean (absolute) amplitude of 5, while the entries have a mean (absolute) value of 1. Yet, it demonstrates the robustness of RMC. Also note that the oversampling has a significant importance in this experiment, as an oversampling of 4 seems to make things harder for RMC: in this case, the RMSE stagnates around  $10^{-4}$ – $10^{-5}$ .

**5.1.3. Noisy low-rank matrix completion with outliers.** In this experiment, we try to tackle the important problem of matrix completion in the presence of *both* (dense) noise and (sparse) outliers. Outliers are defined as previously, while noise is the addition, at each observed entry, of a zero-mean Gaussian random variable with variance  $\sigma_N^2$ .

To have a point of comparison, we compare the results using RMC to the performances that an oracle knowing the row and column space of  $M_0$ , and returning the best matrix using this information, would give. If the entries are perturbed by



(a)  $\mu = \sigma = 0.1$



(b)  $\mu = \sigma = 1$

FIG. 4. Low-rank matrix completion with outliers: Robust low-rank matrix completion of rank-10  $500 \times 500$  matrices observed with an oversampling of 5 and with 5% outliers in the observed entries.

Gaussian noise (without outliers) with variance  $\sigma_N^2$ , the best RMSE is equal (in expectation) to [9]:

$$\text{RMSE}_{\text{Oracle}} = \sigma_N \sqrt{\frac{2nr - r^2}{|\Omega|}}$$

for the low-rank completion of an  $n \times n$  matrix with an  $\ell_2$  method.

Figure 6 depicts the RMSE at termination of RMC with respect to the signal-to-noise (SNR) ratio. Results are the average of three successive experiments. Entries in the matrix  $M$  are such that the matrix has unit-variance Gaussian entries. We thus have  $\text{SNR} = \frac{1}{\sigma_N^2}$ .

We clearly see that—even in the presence of 5% of outliers—the algorithm successfully recovers the original low-rank matrix with an error proportional to the noise level. As long as the noise level is not too high, we have performances very similar to those of the oracle bound. For a high level of noise ( $\text{SNR} < 1$ ), we see that the

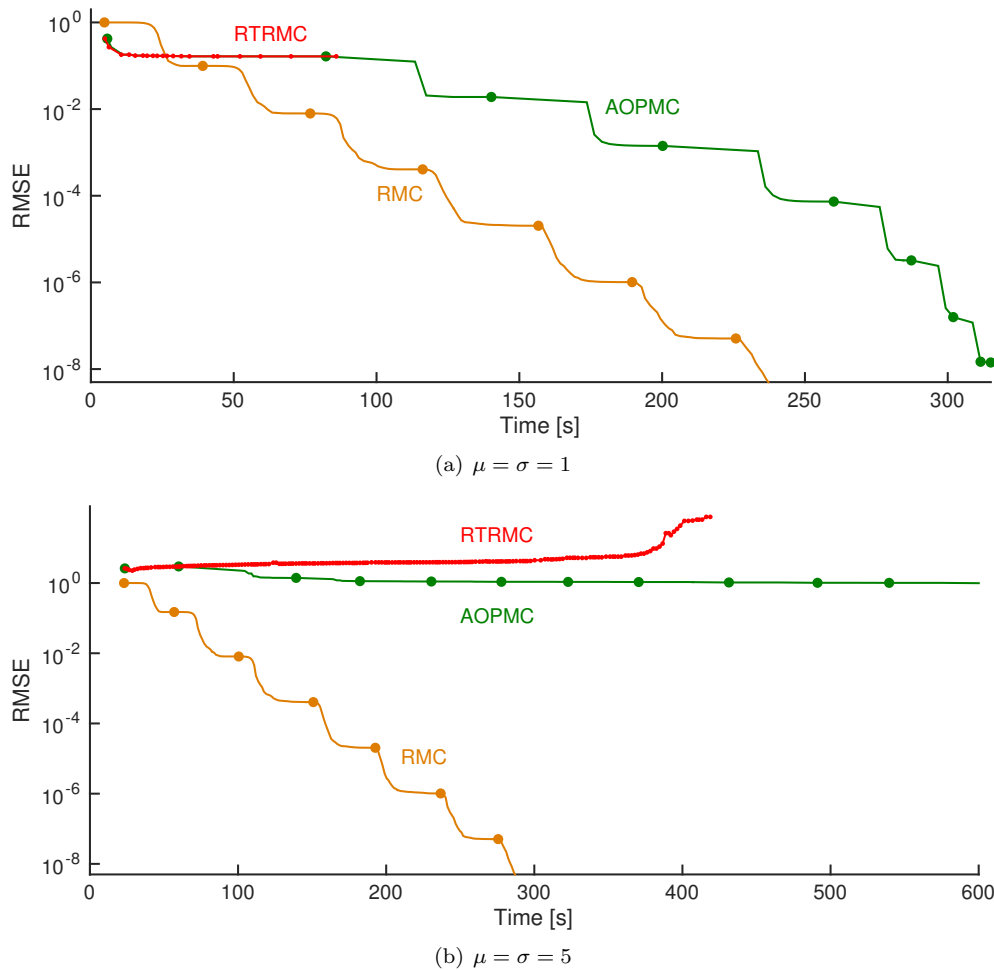


FIG. 5. *Low-rank matrix completion with outliers: Robust low-rank matrix completion on  $50000 \times 50000$  matrices of rank 10 with an oversampling of 5 and 5% outliers in the observed entries.*

algorithm is slightly better than the oracle bound. This can be due to the  $\ell_1$  objective function which helps reduce the effect of the high variance. For a low level of noise, with an SNR greater than  $10^{14}$ , the algorithm begins to have numerical difficulties in driving the RMSE toward zero because the objective function becomes less and less differentiable near the “kinks” of the absolute values. This is the same effect as in the previous experiments where the RMSE begins to stagnate around  $10^{-8}$ – $10^{-6}$  due to the nondifferentiability of the objective function near the solution, which happens to be exactly at the nondifferentiable part of the objective function. A moderate and high level of noise has the effect of “smoothing” the objective function since the solution starts to deviate from the kinks of the  $\ell_1$  norm.

**5.1.4. Influence of the percentage of outliers.** Finally, we study the effect of the percentage and strength of the outliers. As we can expect, for a high number of outliers, we cannot expect RMC to retrieve the original low-rank matrix. But as we will see, as long as the number of outliers is small enough, our algorithm is able to

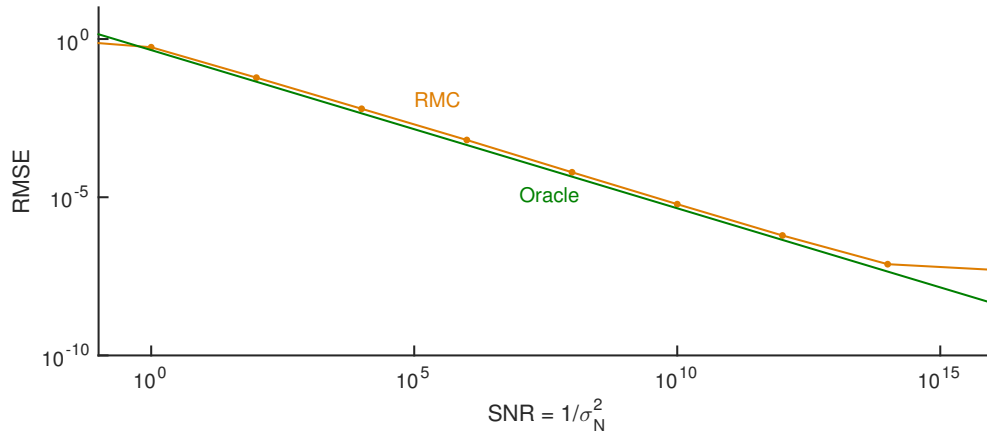


FIG. 6. Noisy low-rank matrix completion with outliers: Evolution of the RMSE with respect to the signal-to-noise ratio  $SNR = \frac{1}{\sigma_N^2}$  on a  $5000 \times 5000$  matrix of rank 10 with an oversampling of 5 and 5% outliers created using  $\mu = 1$  and  $\sigma = 1$ . Noise is the addition of i.i.d. Gaussian variables  $\mathcal{N}(0, \sigma_N^2)$ .

perfectly recover the underlying model. Figure 7 shows how the RMSE at termination evolves with the percentage of outliers added and their strength. This figure is obtained after the average (at each point on the plot) of three experiments, aiming at the completion of a  $5000 \times 5000$  rank-10 matrix observed with an oversampling of 4 (note the triple-log scale).

Three things are worth noting in this figure. First, the abrupt change in the RMSE around 0.5% of outliers is purely numeric. The reason is that the slight increase in the number of outliers allows the method to converge better: the algorithm succeeds in decreasing  $\delta^{(k)}$  one more time (without the CG stalling) and can then decrease the RMSE even more.

Second, we notice a strong increase in the RMSE around 6–8% of outliers. This is quite low but can be explained by the small oversampling ratio used: a low percentage of outliers can reduce the number of healthy entries below the minimum required number. Still, this is interesting, as it shows that as long as the number of outliers is small enough, the recovered matrix stays almost exactly the same as the original unperturbed matrix.

Third, it is interesting to note that the strength of the outliers, from  $\mu = \sigma = 0.1$  to  $\mu = \sigma \approx 10$ , has a negligible impact on the quality of the final solution (remember that entries have values around unity) as long as it remains low enough. This is in complete opposition with the previous experiment on noisy matrix completion, where the RMSE is clearly proportional to the level of the dense Gaussian noise.

**5.2. Application: The Netflix dataset.** One of the most emblematic uses of low-rank matrix completion is in recommender systems and, in particular, in the Netflix Prize [4]. In this section, we propose to test our algorithm RMC on this real-world dataset.

The problem is the following: Netflix, a movie rental company, wants to recommend movies to its users. To do so, they have a limited database of known ratings provided by some users themselves. In practice, this translates into the completion of



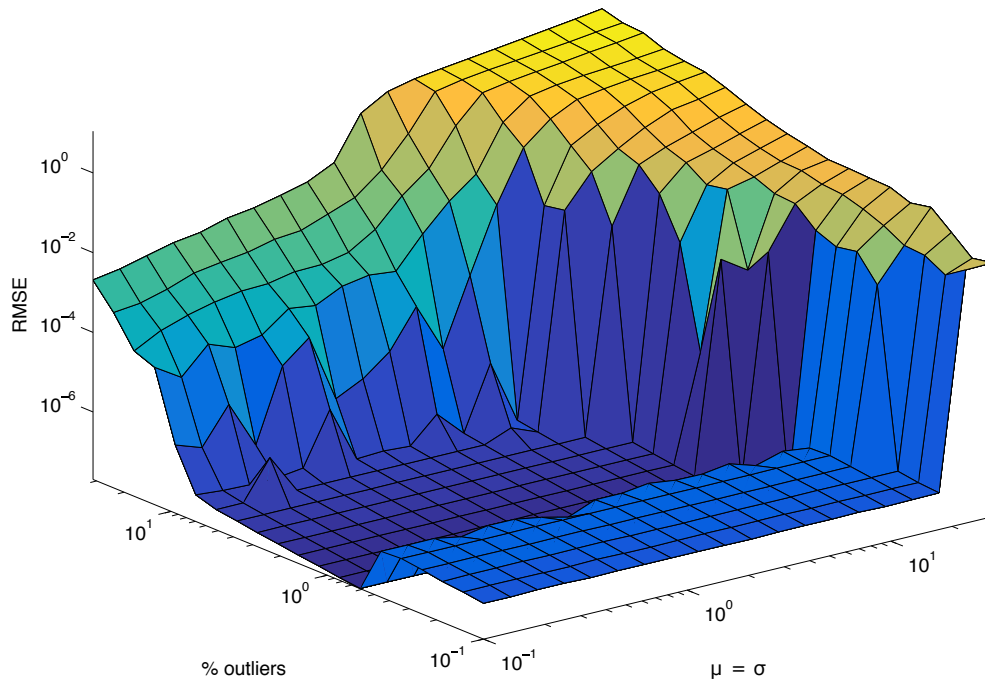


FIG. 7. Evolution of the RMSE with respect to the percentage of outliers and their strength on a  $5000 \times 5000$  matrix of rank 10, observed with an oversampling of 4 and using  $\mu = \sigma$ .

a large  $480189 \times 17770$  matrix, where columns correspond to movies, rows to users, and each entry is the rating of the movie represented by an integer from 1 to 5. To train the algorithm, we have 99072112 entries revealed; that is, approximately 1% of the entries are known. We then test our model on another set of 1408395 entries. Note that all values are shifted toward zero by subtracting the mean of the revealed entries (3.604), since our (regularized) model assumes a mean value of 0.

To assess the results, we use the root mean square criterion, i.e.,

$$\text{RMSE}_{\text{test}} = \sqrt{\frac{\sum_{(i,j) \in \text{TestSet}} (X_{ij} - M_{ij})^2}{|\text{TestSet}|}},$$

on the *test set* (i.e., the unrevealed entries). Returning the mean ratings as a prediction leads to a test RMSE of 1.13, while the winner of the Netflix prize, the BellKor's Pragmatic Chaos algorithm [23], reached an RMSE of 0.8567, using a combination of many different techniques. The authors of [6] performed extensive comparisons between different low-rank matrix completion algorithms (all minimizing the  $\ell_2$  norm on the training set). We can observe that the best result was obtained using LMafit [3], reaching a test RMSE of 0.955.

Two parameters have a significant influence on the results. The rank “dictates” the complexity of the solution; the regularization parameter  $\lambda$  is used to limit overfitting, i.e., to avoid fitting too well the known entries while deviating too much from the mean on the unknown entries. To study the impact of both parameters, we picked a reference point with  $r = 10$  and  $\lambda = 8 \cdot 10^{-4}$ , and we then changed the two parameters around these two values, one at a time. Note that the rank was chosen arbitrarily,

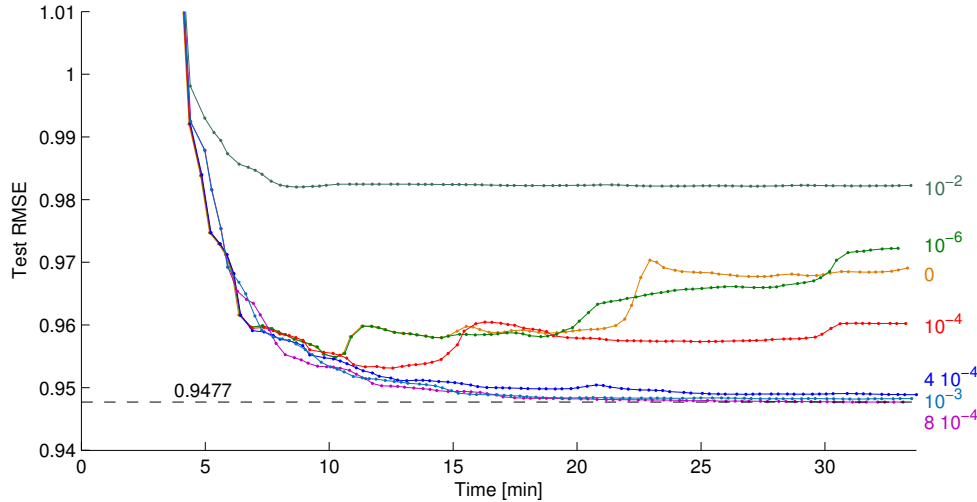


FIG. 8. Evolution of the test RMSE on the Netflix dataset for different values of  $\lambda$  (displayed on the right) using a rank  $r = 10$ . It is clear from this experiment that the regularization  $\lambda$  plays an important role. By tuning it to the right value, we obtain a quite good test RMSE of 0.9477.

while the value of  $\lambda$  is the one that gives the best results for a rank of 10. Also note that, after a few trials, we found that in this context, iteratively decreasing the  $\delta$  parameter was not particularly useful. We then decided to fix it at 1. This value might seem high, but the smooth version is already a quite good approximation of the  $\ell_1$  norm. The number of iterations was limited to 100, and the gradient tolerance was set to  $10^{-8}$ . In practice, most of the runs did not reach a gradient tolerance of  $10^{-8}$  and were interrupted after 100 iterations.

Figure 8 depicts the evolution of the test RMSE for different values of  $\lambda$ . We can easily see that the nonregularized algorithm reaches a reasonable test RMSE but then tends to overfit. Increasing the regularization parameter  $\lambda$  around  $10^{-3}$ – $10^{-4}$  allows us to find a good compromise between training error and overfitting. By increasing it even more, the test RMSE eventually stagnates; i.e., the regularization is so strong that it does not really fit anything except the mean value (see the  $\lambda = 10^{-2}$  curve for instance).

Figure 9 illustrates how the rank plays a significant role in the solution quality. From this plot, it seems that the choice  $r = 10$  was the right one, since other values for  $r$  give higher results. As underlined in [6], choosing a high rank from the beginning does not seem to be the best choice, and rank increasing methods may be worth investigating.

We can conclude from these experiments that our algorithm performs well on the Netflix dataset since it slightly outperforms the low-rank matrix completion algorithms that use the  $\ell_2$  norm. This may be due to the fact that our method is robust to outliers: this can help to reduce overfitting (even with no regularization), hence leading to a better overall model that better fits the test set.

Still, it is known that to reach a lower test RMSE, it is useful to combine this idea of low-rank matrix completion with other techniques. For instance, temporal effects have a large impact, as explained in [4]: movies become more or less popular over

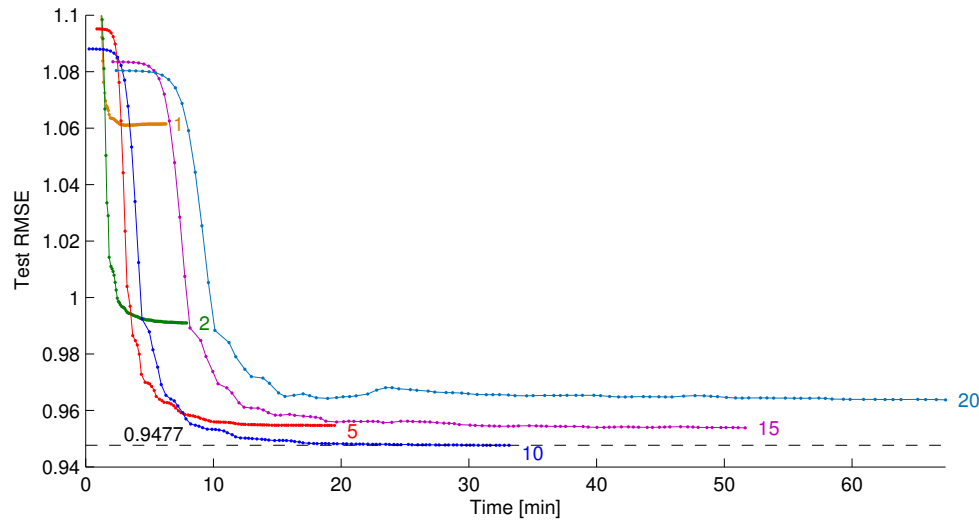


FIG. 9. Evolution of the test RMSE on the Netflix dataset for different values of the rank  $r$  using  $\lambda = 8 \cdot 10^{-4}$ . It seems that the choice  $r = 10$  is the right one, and that increasing the rank leads to overfitting.

time, user tastes and ratings can change over time, and so on. Neighborhood models, where users and movies are aggregated into groups of similar profiles, also help in decreasing the test RMSE [4]. Yet, our algorithm proved itself to be quite efficient and may certainly be used as a good building block for a more complex algorithm. It should also be possible to better fine tune each parameter of the algorithm to find an even better configuration.

**6. Conclusion.** In this paper, we have developed a *robust* low-rank matrix completion algorithm designed to solve the matrix completion problem where a small part of the entries are outliers. Our method can also handle the case where all the entries are corrupted by a small Gaussian noise.

We tackled this problem using Riemannian optimization as well as smoothing techniques to handle the nonsmooth objective function. The resulting algorithm (Algorithm 1) clearly outperforms state-of-the-art methods when outliers have a large amplitude. On the Netflix dataset, it provides a 0.8% improvement on the test RMSE compared to the most accurate (according to [6])  $\ell_2$  low-rank matrix completion method, namely LMaFit.

The code of RMC is available online and can be downloaded from <https://people.stanford.edu/lcambier/rmc>.

**Acknowledgments.** The authors thank Ming Yan and Rahul Mazumder for providing the code of AOPMC and CRMC, respectively, as well as Nicolas Boumal and Yurii Nesterov for their insightful comments.

#### REFERENCES

- [1] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, Princeton, NJ, 2008.

- [2] P.-A. ABSIL AND I. V. OSELEDETS, *Low-rank retractions: A survey and new results*, *Comput. Optim. Appl.*, 62 (2015), pp. 5–29, doi:10.1007/s10589-014-9714-4.
- [3] L. BALZANO, R. NOWAK, AND B. RECHT, *Online identification and tracking of subspaces from highly incomplete information*, in Proceedings of the 2010 48th Annual Allerton Conference on Communication, Control, and Computing, Allerton, IL, IEEE, Piscataway, NJ, 2010, pp. 704–711.
- [4] J. BENNETT AND S. LANNING, *The Netflix Prize*, in Proceedings of the KDD Cup and Workshop, San Jose, CA, 2007, p. 35.
- [5] N. BOUMAL AND P.-A. ABSIL, *RTRMC: A Riemannian trust-region method for low-rank matrix completion*, in Advances in Neural Information Processing Systems 24, Neural Information Processing Systems Foundation, 2011, pp. 406–414.
- [6] N. BOUMAL AND P.-A. ABSIL, *Low-rank matrix completion via preconditioned optimization on the Grassmann manifold*, *Linear Algebra Appl.*, 475 (2015), pp. 200–239, doi:10.1016/j.laa.2015.02.027.
- [7] N. BOUMAL, B. MISHRA, P.-A. ABSIL, AND R. SEPULCHRE, *Manopt, a Matlab toolbox for optimization on manifolds*, *J. Mach. Learn. Res.*, 15 (2014), pp. 1455–1459; toolbox available online at <http://www.manopt.org>.
- [8] E. J. CANDÈS, X. LI, Y. MA, AND J. WRIGHT, *Robust principal component analysis?*, *J. ACM*, 58 (2011), 11, doi:10.1145/1970392.1970395.
- [9] E. J. CANDÈS AND Y. PLAN, *Matrix completion with noise*, *Proc. IEEE*, 98 (2010), pp. 925–936.
- [10] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, *Found. Comput. Math.*, 9 (2009), pp. 717–772, doi:10.1145/2184319.2184343.
- [11] Y. CHEN, A. JALALI, S. SANGHAVI, AND C. CARAMANIS, *Low-rank matrix recovery from errors and erasures*, *IEEE Trans. Inform. Theory*, 59 (2013), pp. 4324–4337, doi:10.1109/TIT.2013.2249572.
- [12] Y. CHEN, H. XU, C. CARAMANIS, AND S. SANGHAVI, *Robust matrix completion with corrupted columns*, in Proceedings of the 28th International Conference on Machine Learning, ACM, New York, 2011, pp. 873–880.
- [13] A. L. CHISTOV AND D. YU. GRIGOR’EV, *Complexity of quantifier elimination in the theory of algebraically closed fields*, in Mathematical Foundations of Computer Science, 1984, Springer, Berlin, 1984, pp. 17–31, doi:10.1007/BFb0030287.
- [14] P. DRINEAS, A. JAVED, M. MAGDON-ISMAIL, G. PANDURANGAN, R. VIRRANKOSKI, AND A. SAVVIDES, *Distance matrix reconstruction from incomplete distance information for sensor network localization*, in Proceedings of the 2006 3rd Annual IEEE Communications Society on Sensor and Ad Hoc Communications and Networks (SECON’06), Vol. 2, 2006, pp. 536–544.
- [15] T. HASTIE, R. MAZUMDER, AND R. TIBSHIRANI, *Matrix Completion and Large-scale SVD Computations*, 2012, [http://web.stanford.edu/~hastie/TALKS/SVD\\_hastie.pdf](http://web.stanford.edu/~hastie/TALKS/SVD_hastie.pdf).
- [16] J. HE, L. BALZANO, AND J. LUI, *Online Robust Subspace Tracking from Partial Information*, preprint, arXiv:1109.3827v2 [cs.IT], 2011.
- [17] J. HE, L. BALZANO, AND A. SZLAM, *Incremental gradient on the Grassmannian for online foreground and background separation in subsampled video*, in Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012, pp. 1568–1575.
- [18] R. KENNEDY, L. BALZANO, S. J. WRIGHT, AND C. J. TAYLOR, *Online algorithms for factorization-based structure from motion*, in Proceedings of the 2014 IEEE Winter Conference on Applications of Computer Vision (WACV), 2014, pp. 37–44.
- [19] R. H. KESHAVAN, A. MONTANARI, AND S. OH, *Low-rank matrix completion with noisy observations: A quantitative comparison*, in Proceedings of the 2009 47th Annual Allerton Conference on Communication, Control, and Computing, Allerton, IL, IEEE, Piscataway, NJ, 2009, pp. 1216–1222.
- [20] O. KLOPP, K. LOUNICI, AND A. B. TSYBAKOV, *Robust Matrix Completion*, preprint, arXiv:1412.8132v1 [math.ST], 2014.
- [21] J. LEE, *Introduction to Smooth Manifolds*, Grad. Texts in Math. 218, Springer-Verlag, New York, 2003, doi:10.1007/978-1-4419-9982-5.
- [22] X. LI, *Compressed sensing and matrix completion with constant proportion of corruptions*, *Constr. Approx.*, 37 (2013), pp. 73–99, doi:10.1007/s00365-012-9176-9.
- [23] NETFLIX, *Netflix Prize Leaderboard*, 2009, <http://www.netflixprize.com/leaderboard>.
- [24] F. NIE, H. HUANG, AND C. DING, *Low-rank matrix recovery via efficient Schatten  $p$ -norm minimization*, in Proceedings of the Twenty-Sixth AAAI Conference on Artificial Intelligence, AAAI, Palo Alto, CA, 2012, pp. 655–661.

- [25] F. NIE, H. WANG, X. CAI, H. HUANG, AND C. DING, *Robust matrix completion via joint Schatten  $p$ -norm and  $lp$ -norm minimization*, in Proceedings of the 2012 IEEE 12th International Conference on Data Mining (ICDM), 2012, pp. 566–574.
- [26] S. OH, A. MONTANARI, AND A. KARBASI, *Sensor network localization from local connectivity: Performance analysis for the mds-map algorithm*, in Proceedings of the 2010 IEEE Information Theory Workshop on Information Theory (ITW), 2010, pp. 1–5.
- [27] Y. PENG, A. GANESH, J. WRIGHT, W. XU, AND Y. MA, *Rasl: Robust alignment by sparse and low-rank decomposition for linearly correlated images*, IEEE Trans. Pattern Anal. Mach. Intell., 34 (2012), pp. 2233–2246.
- [28] A. M.-C. SO AND Y. YE, *Theory of semidefinite programming for sensor network localization*, Math. Program., 109 (2007), pp. 367–384, doi:10.1007/s10107-006-0040-1.
- [29] B. VANDEREYCKEN, *Low-rank matrix completion by Riemannian optimization*, SIAM J. Optim., 23 (2013), pp. 1214–1236, doi:10.1137/110845768.
- [30] B. VANDEREYCKEN, P.-A. ABSIL, AND S. VANDEWALLE, *Embedded geometry of the set of symmetric positive semidefinite matrices of fixed rank*, in Proceedings of the 2009 IEEE/SP 15th Workshop on Statistical Signal Processing, 2009, pp. 389–392.
- [31] M. YAN, Y. YANG, AND S. OSHER, *Exact low-rank matrix completion from sparsely corrupted entries via adaptive outlier pursuit*, J. Sci. Comput., 56 (2013), pp. 433–449, doi:10.1007/s10915-013-9682-3.
- [32] Y. YANG, Y. FENG, AND J. A. K. SUYKENS, *A Nonconvex Relaxation Approach to Robust Matrix Completion*, Internal Report 14-61, ESAT-SISTA, KU Leuven, Leuven, Belgium, 2014.