# Riemannian Gradient Descent Methods for Graph-Regularized Matrix Completion

Shuyu Dong<sup>†</sup>, P.-A. Absil<sup>†</sup>, K. A. Gallivan<sup>‡\*</sup>

<sup>†</sup>ICTEAM Institute, UCLouvain, Belgium <sup>‡</sup>Florida State University, U.S.A.

June 16, 2020

#### Abstract

Low-rank matrix completion is the problem of recovering the missing entries of a data matrix by using the assumption that the true matrix admits a good low-rank approximation. Much attention has been given recently to exploiting correlations between the column/row entities to improve the matrix completion quality. In this paper, we propose preconditioned gradient descent algorithms for solving the low-rank matrix completion problem with graph Laplacian-based regularizers. Experiments on synthetic data show that our approach achieves significant speedup compared to an existing method based on alternating minimization. Experimental results on real world data also show that our methods provide low-rank solutions of similar quality in comparable or less time than the state-of-the-art method.

**Keywords:** matrix completion; graph information; Riemannian optimization; low-rank optimization **MSC:** 90C90, 53B21, 15A83

## 1 Introduction

Low-rank matrix completion arises in applications such as recommender systems, forecasting and imputation of data; see [38] for a recent survey. Given a data matrix with missing entries, the objective of matrix completion can be formulated as the minimization of an error function of a matrix variable with respect to the data matrix restricted to its revealed entries. In various applications, the data matrix either has a rank much lower than its dimensions or can be approximated by a low-rank matrix. As a consequence, restricting the rank of the matrix variable in the matrix completion objective not only corresponds to a reasonable assumption for successful recovery of the data matrix but also reduces the model complexity.

In certain situations, besides the low-rank constraint, it is useful to add a regularization term to the error function, in order to favor other properties related to the true data matrix. This regularization term can often be built from side information, that is, any information associated with row and column indices of the data matrix. Recent efforts in exploiting side information for matrix completion include inductive low-rank matrix completion [53, 23, 58]

<sup>\*</sup>Email:{shuyu.dong, pa.absil}@uclouvain.be, kgallivan@fsu.edu. This work was supported by the Fonds de la Recherche Scientifique – FNRS and the Fonds Wetenschappelijk Onderzoek – Vlaanderen under EOS Project no 30468160. The first author is supported by the FNRS through a FRIA scholarship. Part of this work was performed while the third author was a visiting professor at UCLouvain, funded by the Science and Technology Sector, and with additional support by the Netherlands Organization for Scientific Research.

and graph-regularized matrix completion [61, 26, 59, 42, 55]. In particular, graph-regularized matrix completion involves designing the regularization term using graph Laplacian matrices that encode pairwise similarities between the row and/or column entities (see Section 6.3.1 and Appendix B). Depending on the application, the graph information is available through the connections between the data entities or can be inferred from the data itself.

Rao et al. [42] addressed the task of graph-regularized matrix completion by a low-rank factorization problem of the following form,

$$\underset{(G,H)\in\mathbb{R}^{m\times k}\times\mathbb{R}^{n\times k}}{\operatorname{minimize}} \frac{1}{2} \sum_{(i,j)\in\Omega} \left( (GH^T)_{ij} - M_{ij}^{\star} \right)^2 + \lambda_{\mathrm{r}} \operatorname{Tr} \left( G^T \Theta^{\mathrm{r}} G \right) + \lambda_{\mathrm{c}} \operatorname{Tr} \left( H^T \Theta^{\mathrm{c}} H \right), \quad (1)$$

where  $M^* \in \mathbb{R}^{m \times n}$  is the data matrix to be completed, k is the maximal rank of the low-rank model,  $\Omega$  is the set of revealed entries and  $\lambda_r \geq 0$  and  $\lambda_c \geq 0$  are parameters. The matrices  $\Theta^r := I_m + L^r$  and  $\Theta^c := I_n + L^c$  with given graph Laplacian matrices  $L^r \in \mathbb{R}^{m \times m}$  and  $L^c \in \mathbb{R}^{n \times n}$ . The graph Laplacian matrices  $L^r$  and  $L^c$  incorporate the pairwise correlations or similarities between the columns or rows of the data matrix  $M^*$ . Indeed, observe that the graph Laplacian-based penalty terms in (1) in the form of Tr  $(F^T LF)$  can be written as

$$\operatorname{Tr}(F^{T}LF) = \sum_{i,j} W_{ij} \|F_{i,:} - F_{j,:}\|_{2}^{2},$$
(2)

where W is the graph adjacency matrix such that L = Diag(W1) - W. The right hand-side term above suggests that the graph Laplacian-based penalty term promotes low-rank solutions that show pairwise similarities according to the given graph. Rao et al. [42] related the graph-based regularization term in (1) to a generalized nuclear norm (a weighted atomic norm [8]) and then found a close connection between the matrix factorization model (1) and a convex optimization formulation involving the generalized nuclear norm of the matrix variable  $X = GH^T \in \mathbb{R}^{m \times n}$ . Moreover, they [42, §5] derived an error bound for the generalized nuclear-norm minimization problem, which can be smaller than that of the standard nuclear norm minimization problem if the graph Laplacian matrices are sufficiently informative with respect to the pairwise similarities between the columns/rows of  $M^*$ . In this previous work, an instance (GRALS) of the alternating minimization method is developed for solving the problem (1).

In this paper, we propose to solve the graph-regularized matrix factorization problem (1)by using Riemannian gradient descent and conjugate gradient methods. Our proposed algorithms are motivated by the following consideration. Optimization methods on the matrix product space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  for matrix factorization models have been observed to efficiently provide good quality solutions to matrix completion problems, in spite of the nonconvexity of the cost function. Theoretical support for this observation can be found in [49, 16, 29]. Unlike alternating minimization (e.g. GRALS [42]), both Euclidean gradient descent and our proposed algorithms update the two matrix factors simultaneously, and do not require setting stopping criteria for subproblem solvers as in alternating minimization methods. Furthermore, by exploiting non-Euclidean geometries of the set of low-rank matrices in relation to the matrix product space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , our algorithms use descent directions based on what can be seen as scaled gradients [33, 37] in the matrix product space. Moreover, as in [33], the particular structure of the objective function makes it possible to resort to exact line minimization along the descent direction. We show that the resulting gradient descent algorithms have an iteration complexity bound akin to the Euclidean gradient method (see Theorem 5.5), and that faster convergence behaviors are observed with these proposed algorithms, compared to their counterparts that use the Euclidean geometry (see Section 6).

We test the graph-regularized matrix completion model for matrix recovery tasks on both synthetic and real datasets. We compare our proposed algorithms with a state-ofthe-art method (GRALS [42]), a baseline alternating minimization (AltMin) method and Euclidean gradient descent and conjugate gradient methods. We observe that the proposed algorithms enjoy faster or similar convergence behaviors compared to the state-of-the-art method and faster convergence behavior than the rest of the baseline methods tested. Moreover, the convergence behavior of the proposed algorithms is observed to be more robust against balancing issues that may arise with the asymmetric factorization model in (1), compared to their counterparts that use the Euclidean geometry. For completeness, we also compare empirically the graph-regularized matrix completion model with two low-rank matrix completion models: the matrix factorization model without regularization and the maximum-margin matrix factorization [48, 44] in terms of recovery error. On both synthetic and real datasets, when the graph Laplacian matrices are properly constructed from features of the data matrix (with missing entries), the graph-regularized matrix completion model is found to yield solutions with superior recovery qualities compared to the other two models.

A precursor of this work can be found in the short conference paper [14].

# 2 Related Work and Discussions

Matrix completion models. The graph-regularized matrix completion problem (1) is a generalization of the Maximum-Margin Matrix Factorization (MMMF) problem [48, 44],

$$\underset{(G,H)\in\mathbb{R}^{m\times k}\times\mathbb{R}^{n\times k}}{\text{minimize}} \frac{1}{2} \sum_{(i,j)\in\Omega} \left( (GH^T)_{ij} - M_{ij}^{\star} \right)^2 + \frac{\lambda}{2} \left( \|G\|_F^2 + \|H\|_F^2 \right).$$
(3)

The MMMF problem (3) is related to the nuclear norm-based [5, 43, 7] convex program for low-rank matrix completion [31]

$$\min_{X \in \mathbb{R}^{m \times n}} \frac{1}{2} \sum_{(i,j) \in \Omega} \left( X_{ij} - M_{ij}^{\star} \right)^2 + \lambda \|X\|_*$$

$$\tag{4}$$

via the relation

$$||X||_* = \min_{G,H:GH^T = X} \frac{1}{2} \left( ||G||_F^2 + ||H||_F^2 \right).$$

As shown by Hastie et al. [20], any solution to (3) is also solution to the convex program (4), provided that  $k \ge \operatorname{rank}(M^*)$ . Since the MMMF problem searches for a pair of matrix factors of rank smaller than or equal to k, which is usually much smaller than the matrix dimensions (m and n), its computational cost and memory requirements are significantly reduced compared to the nuclear norm-based convex program (4).

**Optimization methods in related work.** To the best of our knowledge, the stateof-the-art method for graph-regularized matrix completion is the AltMin-type algorithm GRALS [42]. For an alternating minimization method (*e.g.* [42]), one must deal with Sylvester-type equations for solving the subproblem with respect to the low-rank factor G (respectively H) when a graph-based regularization term Tr ( $G^T \Theta^r G$ ) (respectively Tr ( $H^T \Theta^c H$ )) appears in the objective function. Take the fully observed case in [42] for example, the subproblem of (1) with respect to the factor H corresponds to the following Sylvester equation

$$HG^TG + \lambda_c \Theta^c H = M^*G. \tag{5}$$

In the matrix completion scenario, so solving the subproblem of (1) with respect to the factor H corresponds to solving an equation similar to (5), where the constant matrices involved in the equation depend on the positions of the revealed entries in  $\Omega$ . Equivalently,

the subproblem can be rewritten as a linear least-squares problem (in the form of (75) in Appendix A.8) with respect to the vectorization of the factor  $H^T$  of dimension nk. The Hessian operator of this least-squares problem is not block diagonal due to the fact that the original subproblem corresponds to a Sylvester-type equation. Hence the least-squares problem of each alternating minimization step for (1) cannot be decomposed into m or n separate linear systems in dimension k. GRALS [42] approximately solves each of the two least-squares problems by using a linear conjugate gradient (CG) solver.

AltMin-type algorithms have proven to be very efficient in solving bi-convex problems, such as matrix factorization, nonnegative matrix factorization [52, 54], dictionary learning [39, 30], low-rank matrix completion, and in particular have also been proven to converge linearly for the low-rank matrix completion problem (without regularization) [24, 18]. On the other hand, there are heuristic considerations with AltMin-like algorithms in practice. The parameters that control the stopping criteria of the solver for each alternating least squares problem determines the trade-off between the accuracy of the solution and the time efficiency of the AltMin method, but there is no apparent way to set them to achieve the best trade-off once and for all kinds of data. This can be seen in one of our experiments (see Figure 3). GRALS [42], as an instance of the AltMin method with well-tuned parameters and additional stopping criteria in its subproblem solvers, may suffer a significant drop in efficiency when certain properties of the data matrix change: a change of the "scale" of the data matrix, which can be measured by  $||M^*||_F$  for example, changes significantly the performance of GRALS with a fixed set of stopping-criteria parameters.

## **3** Notation, definitions and problem statement

For  $m \in \mathbb{N}^*$ , we denote the set of integers  $\{1, \ldots, m\}$  by  $[\![m]\!]$ . An undirected graph  $\mathcal{G}$ , which is determined by a set of nodes,  $\mathcal{V}$ , a set of (undirected) edges  $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$  and edge weights  $W \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ , is denoted as  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, W)$ . The graph adjacency matrix W is symmetric because  $\mathcal{G}$  is undirected. In addition, we consider adjacency matrices with nonnegative coefficients:

$$W_{i_1,i_2} = W_{i_2,i_1} \begin{cases} > 0 & \text{if } (i_1,i_2) \in \mathcal{E} \\ = 0 & \text{otherwise.} \end{cases}$$
(6)

Throughout this paper, the graph Laplacian matrix of a graph  $\mathcal{G}$ , denoted by L, is defined as

$$L = \operatorname{diag}(W\mathbf{1}) - W. \tag{7}$$

Thus we denote the graph as  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, W)$  or  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, L)$ . The graph Laplacian matrix defined by (6)–(7) is positive semi-definite (e.g. [11, 47]). We denote by  $\Lambda = \text{Diag}(\lambda_1, \ldots, \lambda_{|\mathcal{V}|})$  the diagonal matrix containing the eigenvalues of L in increasing order:  $0 = \lambda_1 \leq \ldots, \leq \lambda_{|\mathcal{V}|}$ . For a matrix  $M \in \mathbb{R}^{m \times n}$ , we model the row index set of M by a graph  $\mathcal{G}^r = (\mathcal{V}^r, \mathcal{E}^r, L^r)$ , where  $\mathcal{V}^r = \llbracket m \rrbracket$ . The superscript r of  $\mathcal{G}^r$  signifies that the graph encodes row-wise correlations. Similarly, the graph that models the column-wise correlations of M is denoted as  $\mathcal{G}^c = (\mathcal{V}^c, \mathcal{E}^c, L^c)$ . For a real-valued symmetric matrix  $\Theta$ , the symbols  $\lambda_{\max}(\Theta)$  and  $\lambda_{\min}(\Theta)$  denote the largest and smallest eigenvalues. The notation  $\Theta \succeq 0$  (respectively  $\Theta \succ 0$ ) signifies that  $\Theta$  is positive semi-definite (respectively positive definite). The largest and smallest singular values of a matrix X are denoted by  $\sigma_{\max}(X)$  and  $\sigma_{\min}(X)$  respectively. The Euclidean inner product and norm for the product space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , denoted as  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|$  respectively, are defined as

$$\langle x, y \rangle = \operatorname{Tr} \left( G_x^T G_y \right) + \operatorname{Tr} \left( H_x^T H_y \right),$$
 (8a)

$$||x|| = \sqrt{\langle x, x \rangle},\tag{8b}$$

for any pair of points  $x = (G_x, H_x), y = (G_y, H_y) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ .

**Problem statement.** The purpose of this paper is to solve (1), which we reformulate as

$$\underset{(G,H)\in\mathbb{R}^{m\times k}\times\mathbb{R}^{n\times k}}{\operatorname{minimize}}\frac{1}{2}\|P_{\Omega}(GH^{T}-M^{\star})\|_{F}^{2}+\frac{\alpha}{2}\left(\operatorname{Tr}\left(G^{T}\Theta^{\mathrm{r}}G\right)+\operatorname{Tr}\left(H^{T}\Theta^{\mathrm{c}}H\right)\right),\tag{9}$$

where  $P_{\Omega}$  is the projection onto the subspace of sparse matrices with nonzeros restricted to the index set  $\Omega$ . The first term in (9) is an equivalent expression of the first term of (1). The second term corresponds to the other terms of (1) with a parameterization that we find more convenient: we set  $\lambda_{\rm r}$  and  $\lambda_{\rm c}$  of (1) by one scalar  $\alpha$ , and  $\Theta^{\rm r} = I_m + \gamma_{\rm r} L^{\rm r}$  and  $\Theta^{\rm c} = I_n + \gamma_{\rm c} L^{\rm c}$ . This allows for more flexible settings of the weight of the Laplacian-based terms over the Frobenius norms. Setting  $\gamma_{\rm r} = 0$  and  $\gamma_{\rm c} = 0$  turns off the graph regularization, leaving us with the MMMF model (3). Setting  $\alpha = 0$  turns off all regularization terms, leading to an unregularized matrix factorization problem

$$\min_{(G,H)\in\mathbb{R}^{m\times k}\times\mathbb{R}^{n\times k}} f_{\Omega}(G,H) := \frac{1}{2} \|P_{\Omega}(GH^T - M^{\star})\|_F^2.$$
(10)

The choice of the rank parameter k is a priori unknown for low-rank matrix approximation problems such as (9). Common approaches include model selection via cross-validation and rank adaptive methods [34, 60]. In this paper, we focus on the setting where k is smaller than or equal to the optimal rank.<sup>1</sup> Observe that (9) is guaranteed to have a solution whenever  $\alpha > 0$  since the objective function is continuous and coercive.

# 4 Optimization on the product space $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$

In this section, we introduce Riemannian gradient descent and conjugate gradient algorithms for problem (9). Since the search space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  is just a vector space, these methods can be interpreted as preconditioned gradient methods. However, we call them "Riemannian" because the preconditioners are inspired from known Riemannian metrics on the set of rank-km-by-n matrices, as we now explain.

In the main problem model (9), the product  $GH^T$  is an  $m \times n$  matrix of rank smaller than or equal to k, and such matrices form the following nonlinear matrix space

$$\mathcal{M}_{\leq k} := \left\{ X \in \mathbb{R}^{m \times n} : \operatorname{rank}(X) \leq k \right\}.$$

In particular, when the regularization parameter  $\alpha$  in (9) reduces to 0, the model (9) can be directly identified with the following optimization problem on  $\mathcal{M}_{\leq k}$  via the matrix factorization model  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k} \mapsto \mathcal{M}_{\leq k} : (G, H) \mapsto X = GH^T$ ,

$$\min_{X \in \mathcal{M}_{\leq k}} \frac{1}{2} \| P_{\Omega}(X - M^{\star}) \|_{F}^{2}.$$
(11)

We refer to the model (11) and (10) as the unregularized matrix completion model, in contrast to the graph-regularized model (9). A recent survey [50] provides advances on optimization methods on the low-rank matrix space  $\mathcal{M}_{\leq k}$ .

In contrast, when the regularization parameter  $\alpha$  in (9) is nonzero, the model (9) does not induce an optimization problem on the Riemannian quotient manifold  $\mathcal{M}_k$  of fixed-rank matrices. This is because the equivalence class of pairs (G, H) that represent the same X is  $\{(GF^T, HF^{-1}) : F \in \mathrm{GL}(k)\}$  and it is readily seen that the graph regularization term (9) is not constant on the equivalence classes. This does not prevent us from drawing inspiration from known Riemannian metrics on the set of rank-k metrices; see next.

<sup>&</sup>lt;sup>1</sup>In real-world applications, it usually suffices to set up a rank value that is orders of magnitude smaller than m and n in order to work with an "underestimated" rank, that is, smaller than or equal to the optimal rank.

**Geometric elements on**  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . The tangent space at a point  $x \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  is the Cartesian product of the tangent spaces of the two element matrix spaces:  $T_x (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}) \simeq \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ .

Given two tangent vectors  $\xi, \eta \in T_x (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$  at x, we consider the following two Riemannian metrics on  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ ,

• The *right-invariant* metric:

$$g_x(\xi,\eta) = \operatorname{Tr}\left(\xi_G^T \eta_G \left(G^T G + \delta I_k\right)^{-1}\right) + \operatorname{Tr}\left(\xi_H^T \eta_H \left(H^T H + \delta I_k\right)^{-1}\right); \quad (12)$$

• The *preconditioned* metric:

$$g_x(\xi,\eta) = \operatorname{Tr}\left(\xi_G^T \eta_G \left(H^T H + \delta I_k\right)\right) + \operatorname{Tr}\left(\xi_H^T \eta_H \left(G^T G + \delta I_k\right)\right),\tag{13}$$

where  $\delta > 0$  is a constant parameter. Both (12) and (13),  $g_x$  are well-defined inner products. The condition  $\delta > 0$  is required for (12) and (13) to remain well defined and positive definite when G or H does not have full column rank. Note that, if we set  $\delta = 0$  and restrict the search space to the product space of full column-rank matrices  $\mathbb{R}^{m \times k}_* \times \mathbb{R}^{n \times k}_*$ , then (12) and (13) reduce to known metrics that induce metrics on the Riemannian quotient manifold  $\mathcal{M}_k$ . Specifically, (12) reduces to the right-invariant metric proposed in [32, 35], and (13) reduces to a metric proposed by Mishra et al. [33], which is specially adapted to the matrix factorization loss function (10). To see this, it suffices to note that the diagonal blocks of the Hessian (see Appendix A.9) of  $f_{\Omega}$  in (10) correspond to the following linear transformation

$$(\xi_G, \xi_H) \quad \mapsto \quad \left( P_{\Omega}(\xi_G H^T) H, P_{\Omega}(G \xi_H^T)^T G \right).$$
(14)

**Definition 4.1.** For a point  $x \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , the gradient of f at x is the unique vector in  $T_x(\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ , denoted as grad f(x), such that

$$g_x(\xi, gradf(x)) = Df(x)[\xi], \forall \xi \in T_x\left(\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}\right).$$
(15)

Based on the metric (12) and Definition 4.1, the gradient  $\operatorname{grad} f(x)$ , denoted as QRIGHT-INV, is

$$\operatorname{grad} f\left(G,H\right) = \left(\partial_{G} f\left(G,H\right) \left(G^{T} G + \delta I_{k}\right), \partial_{H} f\left(G,H\right) \left(H^{T} H + \delta I_{k}\right)\right).$$
(16)

Based on the metric (13) and Definition (4.1), the gradient of  $\operatorname{grad} f(x)$ , denoted as QPRE-CON, is

$$\operatorname{grad} f\left(x\right) = \left(\partial_{G} f\left(G,H\right) \left(H^{T} H + \delta I_{k}\right)^{-1}, \partial_{H} f\left(G,H\right) \left(G^{T} G + \delta I_{k}\right)^{-1}\right).$$
(17)

#### 4.1 Algorithms

In this subsection, we introduce our algorithms and their elements such as computation of the gradient of the objective function of (9) and stepsize selection. We consider two basic algorithms (Algorithms 4.1 and 4.2) for optimization on  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . Computational details of these algorithms are given in Appendices A.1–A.6.

**Initialization.** A widely used initialization method is the so-called spectral initialization (e.g. [27, 28, 49]) to construct the initial low-rank variable  $x^0$ . This consists of computing  $(U_0, S_0, V_0)$  by the k-SVD of the matrix with all the unknown entries set to zero and then defining the initial point  $x^0 := (G^0, H^0)$  as follows,

$$(G^0, H^0) = (U_0 S_0^{1/2}, V_0 S_0^{1/2}).$$
(18)

Algorithm 4.1 Riemannian Gradient Descent (RGD)

Input: Function f: ℝ<sup>m×k</sup> × ℝ<sup>n×k</sup> → ℝ, an initial point x<sup>0</sup> ∈ ℝ<sup>m×k</sup> × ℝ<sup>n×k</sup>, and tolerance parameter ε > 0.
Output: x<sup>t</sup>.
1: t ← 0.
2: Compute gradient: grad f (x<sup>t</sup>).
# See QRIGHTINV (16) or QPRECON (17)

- 3: while  $\|\operatorname{grad} f(x^t)\| > \epsilon$  do
  - For  $\eta^t = -\operatorname{grad} f(x^t)$ , find stepsize  $s_t$  such that  $(s_t, \eta^t)$  satisfy (22) or (23)

# See Algorithm A.4 for (22)

- 5: Update:  $x^{t+1} = x^t + s_t \eta^t$ .
- $6: \quad t \leftarrow t+1.$

7: Compute gradient:  $\operatorname{grad} f(x^t)$ .

8: end while

4:

### Algorithm 4.2 Riemannian Conjugate Gradient (RCG)

**Input:** Function  $f : \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k} \mapsto \mathbb{R}$ , an initial point  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  and  $\epsilon > 0$ . Output:  $x^t$ . 1:  $t \leftarrow 0$ . 2: Compute gradient: grad  $f(x^t)$ ,  $\eta^t = -\text{grad} f(x^t)$ . # See QRIGHTINV (16) or QPRECON (17)3: while  $\|\operatorname{grad} f(x^t)\| > \epsilon$  do Compute the conjuage descent direction  $\eta^t$  by (20). # See Algorithm A.3 4: Find step size  $s_t$  such that  $(s_t, \eta^t)$  satisfy (22) or (23). # See Algorithm A.4 for (22) 5:Update:  $x^{t+1} = x^t + s_t \eta^t$ . 6:  $t \leftarrow t + 1.$ 7: 8: end while

**Gradient and descent directions.** The Euclidean gradient of the objective function of (9) at  $x := (G, H) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  is computed as follows. We have

$$Df(x)[\xi] = \operatorname{Tr}\left(\xi_G^T S H + \xi_H^T S^T G\right) + \alpha \left(\operatorname{Tr}\left(\xi_G^T \Theta^{\mathrm{r}} G\right) + \operatorname{Tr}\left(\xi_H^T \Theta^{\mathrm{c}} H\right)\right),$$

where  $S = P_{\Omega}(GH^T - M^{\star})$ . From the identity

$$Df(x)[\xi] = \operatorname{Tr}\left(\xi_G^T \partial_G f(x)\right) + \operatorname{Tr}\left(\xi_H^T \partial_H f(x)\right)$$

we deduce that the components of the Euclidean gradient  $\nabla f(x)$  are

$$\partial_G f(x) = SH + \alpha \Theta^{\mathbf{r}} G, \tag{19a}$$

$$\partial_H f(x) = S^T G + \alpha \Theta^c H.$$
(19b)

Subsequently, the computation of the Riemannian gradient with respect to the metric (12) (respectively (13)) is based on (16) (respectively (17)) and (19). Algorithms 4.1-4.2 are later referred to as Qrightinv RGD/RCG and Qprecon RGD/RCG respectively. Detailed steps for these computations are given in Algorithm A.1.

In Algorithm 4.1, the descent direction at iteration t is the negative gradient:  $\eta^t = -\text{grad}f(x^t)$ . In Algorithm 4.2, the conjugate descent direction is defined as

$$\eta^{t} = -\operatorname{grad} f\left(x^{t}\right) + \beta_{t} \eta^{t-1} \tag{20}$$

for  $t \geq 1$ , where  $\beta_t$  is determined by a nonlinear CG rule, such as the Fletcher-Reeves rule

$$\beta_t = \frac{g_{x^t} \left( \operatorname{grad} f \left( x^t \right), \operatorname{grad} f \left( x^t \right) \right)}{g_{x^{t-1}} \left( \operatorname{grad} f \left( x^{t-1} \right), \operatorname{grad} f \left( x^{t-1} \right) \right)}, \tag{21}$$

depending on the local geometry of  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . A full description of the computation of the RCG descent directions is given in Appendix A.1; see Algorithm A.3.

**Stepsize selection.** In Algorithms 4.1–4.2, a stepsize  $s_t$  must be selected at each iteration for the update step along a descent direction  $\eta^t \in T_{x^t} (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ . For this purpose, a common approach is to carry out a line search procedure by using backtracking with respect to the Armijo condition,

$$f(x^{t}) - f\left(x^{t} + s\eta^{t}\right) \ge \sigma sg_{x^{t}}\left(-\operatorname{grad} f\left(x^{t}\right), \eta^{t}\right).$$

$$(22)$$

Alternatively, one can estimate the stepsize  $s_t$  via line minimization (e.g. [33, 51]): At a point  $x \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , we compute the stepsize defined as follows,

$$s_t^* = \operatorname*{argmin}_{s \ge 0} f(G + s\eta_G, H + s\eta_H), \tag{23}$$

for a given descent direction  $\eta \in T_x (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ . We use the stepsize (23) for the numerical experiments in this paper. The solution  $s_t^*$  to (23) is obtained by selecting the minimizer from the real positive roots of the derivative of the quartic function of (23), which is a polynomial of degree 3 and can be computed easily. The computational cost of this procedure is of the same order as the computation of the Riemannian gradient (16) or (17). Computational details of (23) are given in Appendix A.4. In Section 5, we will show convergence properties of Algorithm 4.1 using the stepsize (23).

### 4.2 The regularization parameters

In this section, we discuss how to choose suitable values for the parameters of problem (9). Note that the model (9) with  $\alpha > 0$  and at least one strictly positive value for  $\gamma_{\rm r}$  and  $\gamma_{\rm c}$  is referred to as a graph-regularized matrix completion (GRMC) model. When  $\alpha > 0$  and  $(\gamma_{\rm r}, \gamma_{\rm c}) = \mathbf{0}$ , model (9) reduces to a MMMF model (3). When  $(\alpha, \gamma_{\rm r}, \gamma_{\rm c}) = \mathbf{0}$ , model (9) reduces to the unregularized matrix completion model (10), which is referred to as the MC model. Depending on the properties of the data (synthetic and real datasets), and for given graph Laplacian matrices  $L^{\rm r}, L^{\rm c}$ , we have two types of regularization schemes: (i) fixed parameter values and (ii) two-phase regularization scheme.

In the fixed-parameter scheme, we choose the parameter values following a standard cross validation procedure (see, *e.g.*, [25, §5.1.3]). We first generate a collection of parameter settings with random samples drawn from a range of values in  $I_{\alpha} \times I_{\gamma_{\rm r}} \times I_{\gamma_{\rm c}} \subset \mathbb{R}^3_+$  with the uniform distribution in log scale, and then we select the parameter setting that has the best mean validation score (in terms of the RMSE (47) on the validation set). This regularization scheme is used later in Sections 6.2.2 and 6.3.

Algorithm 4.3 Two-phase matrix completion using graph-based regularization (2-phase GRMC)

**Input:** Parameter  $\alpha > 0$ ;  $\gamma_r > 0$  and/or  $\gamma_c > 0$ . Iteration budget T, S > 0 for the two phases. An initial point  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . A tolerance parameter (for Phase 2)  $\epsilon > 0$ . **Output:**  $x^* \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ .

- 1: Initialize with  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  using (18).
- 2: Phase 1:  $x_{\alpha}^* = \text{GRMC Solver}(f_{\alpha}, x^0, \infty)$  for at most T iterations. # See Algorithm 4.1 or 4.2.
- 3: Phase 2: start from  $x_{\alpha}^*$  and find  $x^* = \text{GRMC solver}(f_{\Omega}, x_{\alpha}^*, \epsilon)$  for at most S iterations.

In the two-phase regularization scheme, shown in Algorithm 4.3, we set  $\alpha$  and at least one parameter in  $(\gamma_r, \gamma_c)$  to strictly positive values for Phase 1; for Phase 2, all the regularization parameters are set to zero. In lines 2–3, the GRMC Solver is chosen from one of our proposed algorithms, such as Qprecon RGD (Algorithm 4.1 using the gradient (17)). In Phase 1,  $f_{\alpha}$ denotes the objective function of (9) with the parameter value  $\alpha$ . In Phase 2, the objective function reduces to  $f_{\Omega}$  in (10). The parameters (for Phase 1) in Algorithm 4.3 are chosen in the same way as in the fixed-parameter scheme. This two-phase regularization scheme is designed for sample-efficient exact recovery of low-rank matrices and is used later in Section 6.2.1.

### 5 Convergence analysis

In this section, we analyze the convergence properties of Algorithm 4.1 with step sizes selected by line minimization (23). We conduct the analysis as follows: First, we show that the objective function of (9) is Lipschitz continuously differentiable in the search space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  with respect to the Euclidean geometry. Second, we show that the specially designed non-Euclidean gradient descent directions, defined as QRIGHTINV (16) and QPRE-CON (17), ensure sufficient decrease in the function value provided the step sizes are chosen properly depending on the local geometry at each iterate. We show that the line minimization approach (23) finds such step sizes. Based on these results, we show the convergence behavior of the proposed RGD algorithm based on a generic convergence result given by Boumal et al. [3]. We assume throughout this section that  $\alpha > 0$  in (9). Recall that the inner product  $\langle \cdot, \cdot \rangle$  and the norm  $\|\cdot\|$  throughout this paper are defined in (8) with respect to the Euclidean geometry on  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ .

In the presence of the regularization term, we show that the Euclidean gradient  $\nabla f$  in the sublevel set (with respect to a point  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ ):

$$\mathcal{S}^{0} = \left\{ x \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k} : f(x) \le f(x^{0}) \right\}$$
(24)

is Lipschitz-continuous.

**Lemma 5.1** (Lipschitz-continuous gradient). For a given point  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , there exists a Lipschitz constant  $L_0 > 0$  such that the Euclidean gradient  $\nabla f$  is  $L_0$ -Lipschitz continuous in  $\mathcal{S}^0$  (24): for any  $x, y \in \mathcal{S}^0$ ,

$$f(y) - f(x) \le \langle \nabla f(x), y - x \rangle + \frac{L_0}{2} \|y - x\|^2.$$
 (25)

*Proof.* The objective function of (9) is coercive since  $\alpha > 0$ . Hence  $S^0$  is bounded. Let B be a closed ball that contains  $S^0$ . Since f is  $C^{\infty}$ , it follows that it is Lipschitz continuously differentiable in B. The result follows by a classical argument (see for example [36, Lemma 1.2.3]).

Based on Lemma 5.1, we get the following sufficient decrease property.

**Lemma 5.2.** At any iterate  $x^t = (G^t, H^t)$  produced by Algorithm 4.1 before termination, the following sufficient decrease property holds, provided that the step size s satisfies  $0 < s < 2\Sigma_t/L_0$ , for a strictly positive value  $\Sigma_t > 0$ ,

$$f(x^{t+1}) - f(x^t) \le -C_t(s) \| gradf(x^t) \|^2,$$
(26)

where  $C_t(s) = s(\Sigma_t - \frac{L_0 s}{2}) > 0$ . Under the gradient setting QRIGHTINV (16),

$$\Sigma_t = \min\left(\frac{1}{\delta + \sigma_{\max}^2(G^t)}, \frac{1}{\delta + \sigma_{\max}^2(H^t)}\right),\tag{27}$$

and under the gradient setting QPRECON (17),

$$\Sigma_t = \delta + \min\left(\sigma_{\min}^2(G^t), \sigma_{\min}^2(H^t)\right).$$
(28)

*Proof.* At the t-th iteration in Algorithm 4.1, the descent direction is  $\eta^t = -\text{grad}f(x^t)$ . Let s > 0 denote the step size for producing the next iterate:  $x^{t+1} = x^t + s\eta^t$ . In the gradient setting QPRECON where grad f(x) is defined by (17), the partial differentials are

$$\partial_G f(x) = \eta_G (H^T H + \delta I_k) \text{ and } \partial_H f(x) = \eta_H (G^T G + \delta I_k).$$
 (29)

From Lemma 5.1, we have

$$f(x^{t+1}) - f(x^{t}) \leq \langle \nabla f(x^{t}), x^{t+1} - x^{t} \rangle + \frac{L_{0}}{2} \|x^{t+1} - x^{t}\|^{2},$$
(30)

$$\leq -s \|\eta^t\|^2 \left(\delta + \min\left(\sigma_{\min}^2(G^t), \sigma_{\min}^2(H^t)\right)\right) + \frac{L_0 s^2}{2} \|\eta^t\|^2, \quad (31)$$
  
=  $-C_t(s) \|\operatorname{grad} f(x^t)\|^2. \quad (32)$ 

$$= -C_t(s) \|\operatorname{grad} f\left(x^t\right)\|^2, \tag{32}$$

where  $C_t(s) = s(\Sigma_t - \frac{L_0 s}{2})$  and  $\Sigma_t = \delta + \min(\sigma_{\min}^2(G^t), \sigma_{\min}^2(H^t))$ . The inequality (31) is obtained by using (29) as follows,

$$\langle \nabla f(x^t), x^{t+1} - x^t \rangle = -s \langle \nabla f(x^t), \operatorname{grad} f(x^t) \rangle$$
(33)

$$= -s\left(\operatorname{Tr}\left(\eta_G^T \eta_G(H^T H + \delta I_k)\right) + \operatorname{Tr}\left(\eta_H^T \eta_H(G^T G + \delta I_k)\right)\right) \quad (34)$$

$$\leq -s \left( \delta \|\eta^t\|^2 + \sigma_{\min}^2(H) \|\eta_G\|_F^2 + \sigma_{\min}^2(G) \|\eta_H\|_F^2 \right), \tag{35}$$

where the superscript of the element matrices  $(G, H) = x^t$  and  $(\eta_G, \eta_H) = \eta^t$  are omitted for brevity. Similarly, the same result applies to the gradient setting QPRECON (16), with the quantity  $\Sigma_t$  determined by (27). 

Next, we prove that Algorithm 4.1 with step sizes selected by line minimization (23)ensures sufficient decrease at each iteration.

**Lemma 5.3.** The iterates produced by Algorithm 4.1, with step sizes selected by line minimization (23), satisfy the following sufficient decrease property,

$$f(x^{t+1}) - f(x^t) \le -\left(\Sigma_t^2/2L_0\right) \|gradf(x^t)\|^2.$$
(36)

*Proof.* In Algorithm 4.1, let  $\eta = -\operatorname{grad} f(x^t)$  denote the Riemannian gradient descent direction at iteration  $x^t \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . From Lemma 5.1 and Lemma 5.2, we have

$$f(x^t + s\eta) \le f(x^t) - C_t(s) \|\operatorname{grad} f(x^t)\|^2,$$

for  $s \in [0, 2\Sigma_t/L_0]$ , with  $\Sigma_t$  defined in (28) and (27). One the other hand, let  $\bar{s}$  be the stepsize determined by (23), then by definition, the next iterate  $x^{t+1} = x^t + \bar{s}\eta$  is the minimum of falong the direction  $\eta$ :  $f(x^{t+1}) \leq f(x^t + s\eta)$ , for all  $s \geq 0$ . Hence

$$f(x^{t+1}) \leq \min_{s \in [0, 2\Sigma_t/L_0]} f(x^t + s\eta)$$
(37)

$$\leq \min_{s \in [0, 2\Sigma_t/L_0]} f(x^t) - C_t(s) \| \operatorname{grad} f(x^t) \|^2$$
(38)

$$= f(x^{t}) - \left(\Sigma_{t}^{2}/2L_{0}\right) \| \operatorname{grad} f(x^{t}) \|^{2}.$$
(39)

In both Lemma 5.2 and Lemma 5.3, the sufficient decrease quantity depends on the local information  $\Sigma_t$ . The quantity  $\Sigma_t$  is useful only when it is a strictly positive number. We address this point in Proposition 5.4 for two gradient settings QRIGHTINV (16) and QPRECON (17).

**Proposition 5.4.** Under the same settings as in Lemma 5.2 and 5.3, there exist positive numerical constants  $\Sigma_* > 0$  such that the quantities  $\Sigma_t$  (27) and (28) are lower-bounded,

$$\inf_{t \ge 0} \Sigma_t \ge \Sigma_*. \tag{40}$$

*Proof.* In the case of gradient setting QRIGHTINV (16),

$$\Sigma_t = \min\left(\frac{1}{\delta + \sigma_{\max}^2(G^t)}, \frac{1}{\delta + \sigma_{\max}^2(H^t)}\right).$$

First, we prove that there exists  $D_0 > 0$  such that the norm of the iterate in  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  is bounded:

$$\|x^t\| \le D_0 \tag{41}$$

for all  $t \ge 0$ . It suffices to note that the whole sequence  $(x^t)_{t\ge 0}$  belongs to the sublevel set  $\mathcal{S}^0$  (24) and that f is coercive. Second, for any  $x = (G, H) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , the maximal singular values  $\sigma_{\max}(G)$  and  $\sigma_{\max}(H^t)$  are bounded by the norm as follows,

$$\|x^t\|^2 = \operatorname{Tr}\left(G^T G\right) + \operatorname{Tr}\left(H^T H\right) \ge \sigma_{\max}^2(G) + \sigma_{\max}^2(H).$$
(42)

The result (40) can be deduced by taking the numerical constant  $\Sigma_* := 1/(\delta + D_0^2)$  and combining (41) and (42).

In the case of gradient setting QPRECON (17),

$$\Sigma_t = \delta + \min\left(\sigma_{\min}^2(G^t), \sigma_{\min}^2(H^t)\right) \ge \delta > 0,$$

and the result (40) can be ensured by  $\Sigma_* := \delta$ .

We are now ready to conclude in a manner similar to that in [3, Theorem 2.5]. A minor difference, however, is that the norm in [3] is the Riemannian norm, whereas our search space is a vector space and we use the Euclidean norm (8b). It is possible to use the Riemannian norms induced by (12) and (13) if the iterates stay on a fixed-rank manifold, but in all cases it suffices to use the Euclidean norm in the development of our results.

**Theorem 5.5.** Under the problem statement (9), for a given initial point  $x^0$  and the gradient settings QRIGHTINV (16) and QPRECON (17), the sequence generated by Algorithm 4.1 with the stepsize (23) converges and the decay of the gradient norm satisfies

$$\|gradf\left(x^{N}\right)\| \leq \sqrt{\frac{2L_{0}(f(x^{0}) - f^{\star})}{\Sigma_{\star}N}}$$

$$\tag{43}$$

after N iterations, where  $L_0 > 0$  is the Lipschitz constant mentioned in Lemma 5.1, the numerical constant  $\Sigma_* > 0$  is given in Proposition 5.4 and  $f^*$  is a lower bound<sup>2</sup> of the function value of (9).

<sup>2</sup>One can take  $f^* := 0$  since the objective function of (9) is nonnegative.

*Proof.* The convergence of the sequence  $(x^t)_{t\geq 0}$  is a direct result of the sufficient decrease property (26) in Lemma 5.2 and the boundedness of the sequence of function values  $(f(x^t))_{t\geq 0}$ . See Theorem 2.5 of [3].

Let  $N \ge 1$  denote the number of iterations needed for getting to an iterate  $x^N$  such that  $\|\operatorname{grad} f(x^N)\| \le \epsilon$ , for a tolerance parameter  $\epsilon > 0$ .

Since our algorithm (Algorithm 4.1) does not terminate at  $t \leq N - 1$ , the gradient norms  $\| \operatorname{grad} f(x^t) \| > \epsilon$ , for all  $t \leq N - 1$ . By summing the right hand sides of (36) for  $t = 0, \ldots, N - 1$ , we have

$$f(x^{N}) - f(x^{0}) \leq -\sum_{t=0}^{N-1} (\Sigma_{t}^{2}/2L_{0}) \|\operatorname{grad} f(x^{t})\|^{2}$$
(44)

$$\leq -(\epsilon^2/2L_0)\sum_{t=0}^{N-1}\Sigma_t^2 \tag{45}$$

$$= -(\Sigma_*/2L_0)\,\epsilon^2 N \tag{46}$$

Hence the number of iterations

$$N \le \frac{2L_0(f(x^0) - f(x^N))}{\Sigma_* \epsilon^2} \le \frac{2L_0(f(x^0) - f^*)}{\Sigma_* \epsilon^2}.$$

In other words, the iterate produced by the algorithm after N iterations satisfies

$$\|\operatorname{grad} f(x^N)\| \leq \sqrt{\frac{2L_0(f(x^0) - f^*)}{\Sigma_* N}}.$$

### 6 Numerical Experiments

In this section, we evaluate the performance of the proposed algorithms for solving the graph-regularized matrix completion problem (9). The experimental tests are based on both synthetic and real-world datasets. The synthetic data are generated using a low-rank matrix model with graph information, and the graph Laplacian matrices underlying this graph-related data model are then used in the regularization term of (9). This synthetic experimental setting corresponds to an ideal situation where the graph Laplacian matrices in the regularization term are perfectly conform with the pairwise similarities between the matrix entries. In the tests on real-world data, a graph Laplacian matrix is constructed using basic graph proximity models based on pairwise distances between the rows (or columns) of an initial estimation of the data matrix.

### 6.1 Preliminaries

On both synthetic and real-world data, we compare the time efficiency of our proposed methods with several baseline methods: Euclidean gradient descent on the product space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , alternating minimization and a state-of-the-art method GRALS [42] (also an alternating minimization method). Note that by *time efficiency*, we mean the amount of time that an iterative method takes to arrive at an iterate of a certain recovery accuracy, and this mainly depends on the convergence behavior and the computational cost per-iteration of the method. We also compare the graph-regularized (GRMC) model (9) with two other matrix completion models in terms of matrix recovery errors: (i) unregularized matrix completion (MC) via the factorization model (10) and (ii) the maximum-margin matrix factorization (MMMF) model (3). The following list gives detailed description of the methods involved in the experiments.

• Qprecon RGD and Qprecon RCG correspond to the Riemannian gradient (Algorithm 4.1) and conjugate gradient descent (Algorithm 4.2) using the gradient QPRECON (17). Similarly, Qrightinv RGD and Qrightinv RCG correspond to Algorithm 4.1 and Algorithm 4.2 using the gradient QRIGHTINV (16). By default, the step sizes in all these algorithms are selected via line minimization (23), with a label (linemin). Since we focus on the application of (9) in the low rank setting, we set the rank parameter k by an underestimated value, that is, smaller or equal to the rank of the true hidden matrix, in all experiments. In such case, we set the parameter  $\delta$  in the definition of QPRECON (17) and QRIGHTINV (16) to zero without any numerical issue (*e.g.* having rank deficient factor matrices). We show in Figure 10 (Appendix A.5) that the convergence behavior of the so-tested algorithms is almost the same as their counterparts in the theoretical setting (with a presumably  $\delta > 0$ ) analyzed in Section 5.

For consistency, Euclidean GD and Euclidean CG (see Appendix A.7) stand for the gradient descent and nonlinear conjugate gradient descent algorithms respectively, in which the descent directions (such as the Euclidean gradient) are computed with respect to the Euclidean geometry on  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . The step sizes in all these algorithms are selected via line minimization (23).

- AltMin: An alternating minimization algorithm ((74a)–(74b)), where each of the two graph-regularized least-squares subproblems is solved by the linear CG routine Algorithm A.6. AltMin1 is an instance of AltMin that has accuracy parameter  $\epsilon = 10^{-14}$  and  $n_{\rm CG}^{\rm max} = 500$ . AltMin2 has accuracy parameter  $\epsilon = 10^{-6}$  and  $n_{\rm CG}^{\rm max} = 500$ . Note that the parameter  $n_{\rm CG}^{\rm max}$  can also be set to even larger values: Since each of the two subproblems in the alternating minimization are initialized with the latest iterate (warm-started), the number of iterations required for each subproblem solver to obtain a solution with an accuracy  $\epsilon$  usually does not exceed  $n_{\rm CG}^{\rm max}$  preset here. Hence, the active parameter for controlling the stopping behavior of the subproblem solvers in the experiments is  $\epsilon$ .
- GRALS: An alternating minimization algorithm implemented by Rao et al. [42] that is available on line<sup>3</sup>. GRALS1 denotes the GRALS algorithm with the accuracy parameter  $\epsilon = 10^{-10}$  and  $n_{CG}^{max} = 500$ . GRALS differs from AltMin in that the linear CG solver for the subproblems (74a)–(74b) has an additional stopping criterion<sup>4</sup> compared to AltMin, which could trigger early termination and hence provide inexact solutions to the subproblems (74a)–(74b) under certain circumstances.<sup>5</sup>

To assess the approximation performance for the matrix completion task, we use the root mean-squared-error (RMSE). Given  $M^* \in \mathbb{R}^{m \times n}$  and an index set  $\Omega \subset \llbracket m \rrbracket \times \llbracket n \rrbracket$ , the RMSE of  $X \in \mathbb{R}^{m \times n}$  on  $\Omega$  is defined as

$$\operatorname{RMSE}\left(X;\Omega\right) = \sqrt{\sum_{(i,j)\in\Omega} (X_{ij} - M_{ij}^{\star})^2 / |\Omega|}.$$
(47)

<sup>&</sup>lt;sup>3</sup>Link: https://github.com/rofuyu/exp-grmf-nips15.

<sup>&</sup>lt;sup>4</sup>A stopping criterion that restricts the subproblem update to a region of radius depending on the norm of the partial gradients of f.

<sup>&</sup>lt;sup>5</sup>This feature is one of reasons that make the convergence behavior of GRALS different from that of AltMin, and in several applications, faster than the latter.

All numerical experiments are performed on a workstation with 8-core Intel Core i7-4790 CPUs and 32GB of memory running Ubuntu 16.04 and MATLAB R2015a. The source code is available on https://gitlab.com/shuyudong.x11/grmc.

#### 6.2 Synthetic Data

We generate synthetic data with the following low-rank matrix model, which is a generalization of the model in [42, §5.1]. Let  $\mathcal{G}^{\mathrm{r}} := (\mathcal{V}^{\mathrm{r}}, \mathcal{E}^{\mathrm{r}}, L^{\mathrm{r}})$  and  $\mathcal{G}^{\mathrm{c}} := (\mathcal{V}^{\mathrm{c}}, \mathcal{E}^{\mathrm{c}}, L^{\mathrm{c}})$  be the graphs modeling the row-wise and column-wise similarities of  $M^{\star}$  and let  $(U^{\mathrm{r}}, \Lambda^{\mathrm{r}})$  (respectively  $(U^{\mathrm{c}}, \Lambda^{\mathrm{c}})$ ) denote the pair of matrices containing the eigenvectors and associated eigenvalues of  $L^{\mathrm{r}}$  (respectively  $L^{\mathrm{c}}$ ). A low-rank data matrix  $X^{\star}$  is generated as follows,

$$Z^{\star} = F^{\star}Q^{\star T}, \tag{48a}$$

$$X^{\star} = A^{\mathrm{r}} Z^{\star} (A^{\mathrm{c}})^{T}, \qquad (48\mathrm{b})$$

where  $(F^*, Q^*) \in \mathbb{R}^{m \times r^*} \times \mathbb{R}^{n \times r^*}$  are composed of columns that are *i.i.d.* Gaussian vectors and the matrices  $A^{\mathbf{r}} \in \mathbb{R}^{m \times m}$  and  $A^{\mathbf{c}} \in \mathbb{R}^{n \times n}$  are defined below with respect to a function  $g : \mathbb{R} \mapsto \mathbb{R}$  acting element-wisely on a diagonal matrix:

$$A^{\mathbf{r}} = U^{\mathbf{r}}g(\Lambda^{\mathbf{r}}), A^{\mathbf{c}} = U^{\mathbf{c}}g(\Lambda^{\mathbf{c}}).$$
(49)

More precisely,  $g(\Lambda) = \text{Diag}(g(\lambda_1), \dots, g(\lambda_m))$  for any diagonal matrix  $\Lambda$ . The function g in (49) enables one to control the way in which the graph information in  $L^r$  and  $L^c$  transforms the low-rank random Gaussian matrix  $Z^*$ . In the literature of graph signal processing [45], the function g is referred to as a graph spectral filter, which is a graph analogue of filters in signal processing. In our experiments, the function g is

$$g(\lambda) = \begin{cases} \lambda^{-p} & \text{if } \lambda > 0, \\ 0 & \lambda = 0, \end{cases}$$
(50)

for  $p \geq 1$ . The spectral model (50) is a typical example of functions that are monotonically non-increasing over  $\mathbb{R}^*_+$  and that have the effect of low-pass filters [45] in the graph spectral domain [46]. Other examples include (i) the Tikhonov filter (*e.g.* [1])  $g_{\gamma}(\lambda) = 1/\sqrt{1+\gamma\lambda}$ , and (ii) the diffusion operator [12, 13, 56]  $g_{\tau}(\lambda) = e^{-\tau\lambda}$ .

**Remark 6.1.** In order for model (48) to cover the graph-agnostic setting as a special case, we define by convention that  $A^{\rm r} = I_m$  when  $L^{\rm r} = \mathbf{0}$  and  $A^{\rm c} = I_m$  when  $L^{\rm c} = \mathbf{0}$ .

The model (48) is of particular interest for experiments on synthetic data because it models a wide range of real data matrices whose entries present pairwise similarities: Due to the non-increasing nature of the function g on  $(0, \infty)$  in (49), the transformations in (48b) with the matrices  $A^{\rm r}$  and  $A^{\rm c}$  return a data matrix  $X^*$  such that the graph-based regularization terms of (9) are reduced compared to that before the transformation (see Appendix B for details). This translates to the observation that the entries of  $X^*$  present pairwise similarities that agree with the graph  $(\mathcal{V}^{\rm r}, L^{\rm r})$  and/or  $(\mathcal{V}^{\rm c}, L^{\rm c})$ , unlike the structureless entries in  $Z^*$  (48a). Figure 1 shows the difference between  $Z^*$  and  $X^* := A^{\rm r}Z^*$  regarding this property.

#### 6.2.1 Matrix completion from noiseless observations

In this subsection, the ground-truth data matrix  $M^*$  is generated by (48) for  $r^* \ll \min(m, n)$ and is partially observed without any noise. The index of the revealed entries are *i.i.d.* sampled according to the Bernoulli model

$$(i, j) \in \Omega$$
 with probability  $\rho$ , for any  $(i, j) \in \llbracket m \rrbracket \times \llbracket n \rrbracket$ . (51)



Figure 1: A data matrix Z from the Gaussian random model (48a) and  $X = A^r Z$  from the graph-based model (48b). The Laplacian matrix  $L^r$  involved in (48b) is generated with the prototypical graph model Community using GSPbox [40]. Comparison in the data entries and in the graph spectral domain. Top: A randomly chosen column (a):  $y = Z_{(j)}$  and (b):  $y = X_{(j)}$ . Bottom: Average amplitude of graph Fourier coefficients (c):  $\tilde{\mathbb{E}}|\widehat{Z_{(j)}}(\lambda_l)|$  and (d):  $\tilde{\mathbb{E}}|\widehat{X_{(j)}}(\lambda_l)|$  from low (small eigenvalue  $\lambda_l(L^r)$ ) to high graph vertex frequencies.

For simplicity, we let  $L^{c} = 0$  such that  $A^{c} = I_{n}$  (see Remark 6.1). Hence the graph information is incorporated in  $M^{*}$  row-wisely by  $A^{r}$  with respect to (49). For this purpose, a graph Laplacian matrix  $L^{r}$  is generated with the prototypical graph model Community using the GSPbox [40]. The function g in this model is (50) with p = 2.

For the matrix completion model (9), we set the rank parameter by  $k := \operatorname{rank}(M^*)$ . Note that in this case,  $M^*$  belongs to  $\mathcal{M}_{\leq k}$ , and any point  $(G^*, H^*) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  such that  $G^*H^{*T} = M^*$  exactly recovers the hidden matrix  $M^*$ . We refer to the search for such a point  $(G^*, H^*)$  as exact recovery of the data matrix. In the literature of matrix completion, exact recovery of a low-rank matrix  $M^*$  by a factorization model such as (10) is possible under conditions on the extent of incoherence [6] of the singular subspaces of  $M^*$  and the observation model  $\Omega$ . Specifically, several sample complexity lower-bounds for  $\rho \approx |\Omega|/mn$ are proved with both regularized ([49, 16]) and unregularized (implicitly regularized [29]) matrix factorization models.

In the experiments of this subsection, we carry out tests for recovering the hidden matrix  $M^*$  with our proposed two-phase (2-phase GRMC) regularization scheme (Algorithm 4.3). Note that this 2-phase regularization scheme is specially adapted to the exact recovery of the hidden matrix  $M^*$  since it disables the regularization terms in its last phase (avoiding any bias in the solution). The unregularized matrix completion (MC) model (10), which corresponds to the special setting of (9) for  $(\alpha, \gamma^r, \gamma^c) = \mathbf{0}$ , is also tested. The label "MC (GRALS)" in Figure 2(a) corresponds to the result of unregularized matrix completion using **GRALS**, which reduces to a simple "ALS" algorithm since all regularization parameters are set to zero for the (unregularized) MC model.



Figure 2: Matrix completion from noiseless observations.  $M^*$  is generated with non-trivial graph information with the model (50) and is partially observed without any noise. The rank parameter  $k := \operatorname{rank}(M^*)$ . (a): Percentage of successful recoveries under various sampling rates. The solutions are given by the two different matrix completion models (MC and GRMC). Matrix size m = 500, n = 600, rank  $r^* = 12$ . (b): Results per iteration by 2-phase GRMC (Algorithm 4.3) at sampling rate  $|\Omega|/mn = 10.0\%$ : matrix size m = 800, n = 900, rank  $r^* = 12$ .

First, we compare empirically the sample complexities of (i) the unregularized matrix completion model (MC) and (ii) the graph-regularized matrix completion model (9) through the 2-phase GRMC scheme described above. Under the experimental settings described in the beginning of this Section, for m = 500, n = 600 and  $r^* = 12$ , we carry out repeated tests at various sampling rates  $|\Omega|/mn$  ranging from 5% to 28%. At each sampling rate, we compute the percentage of successful recoveries among  $N_{\text{tests}} = 20$  repeated tests. Each test is counted as successful if the RMSE (47) on test entries is smaller than  $10^{-12}$ .<sup>6</sup> In particular, at each sampling rate, the parameters  $(\alpha, \gamma_r)$  in the GRMC model (9) (for Phase 1 of Algorithm 4.3) are selected with respect to the test RMSE of the final solution, among NCONFIGS = 5 randomly generated parameter configurations (see the paragraph of the fixed-parameter scheme of Section 4.2 for details). Here the configurations for  $(\alpha, \gamma_r)$  are generated with the uniform distribution (in the log scale) in the 2-dimensional box  $[10^{-4}, 1] \times [10^{-2}, 5]$ .<sup>7</sup>

As shown in the 2-phase GRMC scheme (Algorithm 4.3), the whole algorithm is stopped by either the accuracy parameter  $\epsilon$  (see Algorithms 4.1, 4.2), when the iterate becomes an  $\epsilon$ -stationary point or by the iteration budget parameter S, which is tuned to a sufficiently large value for both successful and unsuccessful recovery scenarios. Experimental results are shown in Figure 2(a): These results show empirically that the 2-phase GRMC method has a lower sample complexity than unregularized matrix factorization.

Second, we compare the time efficiency of the proposed algorithms with their counterparts under the Euclidean geometry. Figure 2(b) shows results per iteration under a sampling rate that is sufficiently large for 2-phase GRMC.

In particular, when the sampling rate  $\rho$  is sufficiently large, it is possible to exactly recover the hidden matrix  $M^*$  without any regularization (see Figure 2(a) at sampling rates

<sup>&</sup>lt;sup>6</sup>This is an attainable accuracy level in the exact recovery scenario, based on preliminary tests.

<sup>&</sup>lt;sup>7</sup>In the exact recovery scenario, the Phase 1 of our 2-phase algorithms does not need very fine-tuned parameters and the 2-dimensional box was also already narrowed after preliminary tests. A selection from 5 parameter configurations was enough to get the improvements shown in Fig.2(a).

larger than 15%). Therefore, we test our algorithms for this special case, with the problem parameter  $\alpha$  set to zero in (9). In this special regime, we compare the time efficiency of our proposed algorithms with the several other methods in two different settings for the initialization point  $x^0 \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ : In the first test, each method is initialized at a point  $x^0 = (G_0, H_0)$  given by (18). In this case, the two factors  $G_0$  and  $H_0$  are *balanced*, in the sense that their matrix norms are equal. In the second test, we test the same methods with an unbalanced initial point

$$y^0 = (\lambda G_0, H_0/\lambda), \tag{52}$$

for  $\lambda = 5$ . The comparative results are given in Figure 3.



Figure 3: Results per iteration. Experimental settings:  $m = 1000, n = 900, r = r^* = 10, |\Omega|/mn = 20.0\%$ .  $M^*$  is generated with the model (50) and is partially observed without any noise. (a): Each method is initialized at  $x^0$  by (18). (b): Each method is initialized at  $y^0$  by (52).

From the results in Figure 2 and Figure 3, we have the following observations:

- Our algorithms (Qprecon RGD, RCG) are faster than their Euclidean geometry-based counterparts (Euclidean GD, CG) in every experimental setting.
- Our algorithms are faster than the baseline alternating minimization methods AltMin1, AltMin2.
- Qprecon RCG is faster than GRALS1 and this comparison becomes much more evident when the initialization point is unbalanced than when it is balanced. Similarly, Qprecon RGD is as fast as GRALS1 in the balanced initialization setting and much faster than the latter in the unbalanced case.
- In relation to the remark above, the baseline methods Euclidean GD, Euclidean CG and GRALS1 are significantly slower when the initialization point is unbalanced.

### 6.2.2 Matrix completion from noisy observations

In this subsection, we assume that the partially observed data matrix  $M^*$  is composed of noisy observations from a low-rank matrix  $X^*$ ,

$$M^{\star} = X^{\star} + E, \tag{53}$$

where  $E_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_N^2)$  for all  $(i, j) \in \llbracket m \rrbracket \times \llbracket n \rrbracket$  and  $X^* \in \mathbb{R}^{m \times n}$  is generated using the model (48) with rank  $r^* \ll \min(m, n)$ . For simplicity, we let  $L^c = 0$  and only incorporate row-wise similarities in  $X^*$  through  $L^r \in \mathbb{R}^{m \times m}$  with respect to (49). For this purpose, a graph Laplacian matrix  $L^r$  is generated with the prototypical graph model Community using the GSPbox [40]. The function g in this model is (50) with p = 2. Figure 4 shows the singular values of  $M^*$  generated with (53), where the true low-rank matrix  $X^*$  is generated with (48), and the noise level of E is determined by a signal-to-ratio parameter SNR = 20.



Figure 4: Singular values of  $M^*$  generated with (48) and (53), where the noise level is set by a signal-to-ratio parameter SNR = 20.

For the matrix completion problem (9), the rank parameter k is set to be smaller than  $\operatorname{rank}(X^*)$  and its values will be specified later. We test our algorithms, the baseline algorithms and the state-of-the-art algorithm GRALS [42] for solving the problem (9) with fixed parameter values  $\alpha \geq 0$  and  $\gamma_r \geq 0$ .

First, we compare the recovery quality of the solutions to the three types of matrix completion models, that is, the unregularized matrix completion (MC), the graph-regularized matrix completion (GRMC) model (9) and the maximum-margin matrix completion (MMMF) in (3). For the unregularized problem model (MC), the parameter setting is  $\alpha = 0$ . For the MMMF model ( $\alpha > 0$  and  $\gamma_r = 0$ ), the parameter setting is selected among a collection of  $N_{\rm HP}$  randomly generated values in a reasonable range,  $(\alpha^i)_{i=1,\dots,N_{\rm HP}}$ . For the GRMC model  $(\alpha > 0 \text{ and } \gamma_{\rm r} > 0)$ , the parameter setting is selected among a collection of NCONFIGS = 10 uniformly distributed (in log-scale) values in a 2-dimensional box  $[10^{-4}, 1] \times [10^{-2}, 5]$ . The criterion for the parameter selection is the test RMSE (47). At each sampling rate and once the parameters are selected for MMMF and GRMC, the recovery score of each of the three matrix completion models, using a given algorithm (ours as well as the baseline methods) corresponds to the average score after  $N_{\text{tests}}$  training instances based on  $(M^{\star}, \Omega_s)_{s=1,\dots,N_{\text{tests}}}$ , where the observation sets  $\Omega_s$  are generated according to (51) with the given (fixed) sampling rate. All methods tested are stopped by either the accuracy parameter  $\epsilon$  (see Algorithms 4.1, 4.2), when the iterate becomes an  $\epsilon$ -stationary point or by an iteration budget parameter, which is set to a sufficiently large value for all methods. Figure 5 shows the recovery scores of each of the three problem models under different sampling rates. From Figure 5, we can see that at various sampling rates, the GRMC model (9) provide solutions with superior recovery qualities than the other two graph-agnostic models. Naturally enough, the improvement on recovery qualities via GRMC is significant at small sampling rates.

Second, we compare the time efficiency of the proposed algorithms with the baseline methods. The methods are tested in two slightly different experimental settings. Based on the data generation method described in the beginning of this subsection, the data matrix



Figure 5: Average test RMSE of the three matrix completion models from noisy observations.  $M^*$  is generated with non-trivial graph information with the model (50) and is partially observed with additive noise (SNR = 20). The rank parameter k = 10. Recovery score of solutions to MC, MMMF and GRMC at various sampling rates. Matrix size: m = 500, n = 600, rank  $r^* = 12$ . The sampling rate  $|\Omega|/mn$  ranges from 1.0% to 18.0%.

in each of these two experiments is rescaled with respect to a given constant value: We set the constant scalar such that  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 1$  in the first setting and  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 10^{-3}$  in the second one. Figure 6 shows the time efficiency of the tested methods in these two experiments in terms of the RMSE score per iteration. Note that the tests for producing Figure 5 are conducted under the data setting  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 1$ , without loss of generality. In particular, for the second data setting, with  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 10^{-3}$ , we show in Figure 7 the recovery qualities of the iterates under the unregularized (MC) and the graph-regularized (GRMC) models, for a relatively low sampling rate. Given the graph  $L^{r}$  underlying the synthetic data model (48), the GRMC model corresponds to one randomly generated set of parameters ( $\alpha > 0, \gamma_{r} > 0$ ), where  $\alpha$  is randomly generated in the range  $(10^{-6}, 10^{-3})$  and  $\gamma_{r}$  randomly generated in the range  $(10^{-2}, 5)$ . We can see that for all the tested methods, the recovery qualities of the iterates under the GRMC model outperforms those of the unregularized matrix completion model.

From the results in Figure 6 and Figure 7, we have the following observations:

- Our algorithms (Qprecon RGD, RCG) are faster than their counterparts under the Euclidean geometry (Euclidean GD, CG).
- Our algorithms are faster than the baseline alternating minimization methods AltMin1, AltMin2.
- Qprecon RCG is either faster than GRALS1 or as fast as the latter in various settings.
- In relation to the previous remark. The time efficiency of GRALS changes significantly when there is a simple change in the scale of the data matrix, as shown in Figure 6, since it has an additional stopping criterion that restricts the search of the solution to the least-squares subproblem (74a) (resp. (74b)) to a region of radius  $\|\partial_G f(G^{t-1}, H^{t-1})\|$ (respectively  $\|\partial_H f(G^t, H^t)\|$ ). This restricted-region criterion depends however, both on properties of the data matrix (such as the scale of the data) and the iterate.
- As illustrated in Figure 5, the recovery performances of all tested methods for GRMC are better and more stable than those for (unregularized) MC, when the sampling rate  $|\Omega|/mn$  is insufficient; see Figure 7(a)–(b).



Figure 6: Test RMSE per-iteration on sythetic data. The data matrix  $M^{\star}$  is generated via (48) and (53). Matrix size: m = 1000, n = 900, rank  $r^{\star} = 12$ . The rank parameter k = 8. Subplots (a-b): The data matrix  $M^{\star}$  is rescaled by a scalar constant such that  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 1$ ; (a): the sampling rate  $|\Omega|/mn = 11.5\%$ , (b):  $|\Omega|/mn = 18.0\%$ . Subplots (c-d): The data matrix  $M^{\star}$  is rescaled by a scalar constant such that  $\mathbb{E}\left[|M_{ij}^{\star}|\right] = 10^{-3}$ ; (c): the sampling rate  $|\Omega|/mn = 18.0\%$ .



Figure 7: Test RMSE per iteration on synthetic data. The data matrix  $M^*$  is rescaled by a scalar constant such that  $\mathbb{E}\left[|M_{ij}^*|\right] = 10^{-3}$ . Matrix size: m = 1000, n = 900, rank  $r^* = 12$ . The rank parameter k = 8. (a)–(c): the sampling rate  $|\Omega|/mn$  is 3.5%, 5% and 10% respectively. The dashed lines with a label "(MC)" corresponds to methods for solving the (unregularized) MC problem.

### 6.3 Real Data

In this subsection, we conduct experiments on real-world datasets. An essential difference between these experiments and experiments on synthetic data is that there is no reference graph associated with the data matrices in a real-world application. Since the data matrix in a real-world application often present pairwise similarities between its entries, we build the graphs  $L^{\rm r}$  and  $L^{\rm c}$  based on the given data before the graph-regularized matrix completion task. Subsequently, we conduct tests with the graph-regularized and graph-agnostic matrix completion models using all the methods involved, and compare their time efficiency. The real-world data used for these tests are from the PeMS Traffic occupancy and MovieLens datasets.

#### 6.3.1 Methodology for graph construction

In the existing work on matrix completion using graph-based regularization, there are two main approaches to constructing the graph Laplacian matrices: (i) build the graph Laplacian matrix from the data  $M^* \in \mathbb{R}^{m \times n}$  itself by using a certain graph node proximity model, *e.g.*, [26], and (ii) build a similarity graph  $L^r$  (and/or  $L^c$ ) from side information [42, 55], that is, information related to the entities of the row (and/or column) indices of  $M^*$ .

In our experiments, we adopt the first approach. Note that in [26], the computation of the graph proximity parameters is based on pairwise distances using only the revealed entries in  $M^*$ . In contrast, we compute the graph proximity parameters based on a low-rank approximation of the partially revealed matrix. More precisely, we propose to use a rank-rapproximation of  $M_0$  as the features for constructing the graph. Let  $(U_0, S_0, V_0)$  denote the r-SVD of the zero-filled matrix  $M_0 := P_{\Omega}(M^*) \in \mathbb{R}^{m \times n}$  and let  $\widetilde{M}_0 := U_0 S_0 V_0^T$ . Next, the computation of the graph edge weight parameters based on the given matrix  $M := \widetilde{M}_0$  can be realized by using various node proximity methods such as K-Nearest Neighbors (K-NN) and  $\varepsilon$ -graph models [9, 2, 21, 10], which boils down to computing a certain distance matrix between the rows (respectively columns) of M. Let  $Z^{\mathrm{r}}(M) \in \mathbb{R}^{m \times m}$  denote the row-wise distance matrix of M defined as follows,

$$Z_{ij}(M) = d(M_{i,:}, M_{j,:}), \text{ for } i, j \in [\![m]\!],$$
(54)

where  $d : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}_+$  is a distance on the *n*-dimensional vector space. Subsequently, we build a Gaussian  $\varepsilon$ -graph by computing the node proximity weights as follows

$$[W_{\varepsilon}(M)]_{ij} = \exp\left(-Z_{ij}(M)/\varepsilon^2\right), \text{ for } i, j \in [\![m]\!],$$
(55)

where  $\varepsilon \in \mathbb{R}$  is a hyperparameter of the graph model. Furthermore, a sparse graph adjacency matrix is more preferable than a dense in a computational point of view, as the per-iteration cost for computing the gradient (as well as the function value) in (19) depends partly on  $nnz(L^r)$  and  $nnz(L^c)$ , hence the sparsity of the row-wise (resp. column-wise) graphs. For simplicity, we sparsify the graph adjacency matrix defined in (55) by the following thresholding operation

$$[W_{\varepsilon,\sigma}(M)]_{ij} = \mathbf{1}_{\geq \sigma} \left( \exp\left(-Z_{ij}(M)/\varepsilon^2\right) \right), \text{ for } i, j \in [\![m]\!],$$
(56)

where  $\mathbf{1}_{\geq \sigma}$  is the hard threshold function  $\mathbf{1}_{\geq \sigma}(z) = \begin{cases} z & \text{if } z \geq \sigma \\ 0 & \text{otherwise.} \end{cases}$ 

In the graph model (56), the parameter  $\varepsilon$  is tuned according to the variance of  $(Z_{ij})_{i,j=1,...,m}$ . In the following experiments, we set  $\varepsilon := \operatorname{Var}((Z_{ij})_{ij})/5$  and find that this setting gives satisfactory improvements on the final recovery performances. The parameter  $\sigma$  is chosen according to a preset sparsity level  $\mathfrak{s} \ll 1$  for the edge set associated with  $W_{\epsilon,\sigma}$  such that  $|\mathcal{E}(W_{\epsilon,\sigma})|/m^2 \leq \mathfrak{s}$ . We set  $\mathfrak{s} := 8\%$  and find that it is an appropriate trade-off between the amount of similarity graph edges and the additional computational cost required by matrix multiplications with the graph Laplacian.

### 6.3.2 The Traffic Data

The PeMS Traffic occupancy data<sup>8</sup> is a matrix with dimensions  $963 \times 10560$  containing traffic occupancy rates (between 0 and 1) recorded across time by m = 963 sensors placed along different car lanes of the San Francisco Bay area freeways. The recordings are sampled every 10 minutes covering a period of 15 months. The column index set corresponds to the time domain and the row index set corresponds to geographical points (sensors), which are referred to as the spatial domain. Unlike the case with data from social networks or any other kind with useful meta-data, there is no straightforward way to find any side information for the Traffic dataset that may help constructing a spatial-domain graph. Hence we construct a sparse row-wise similarity graph with the Gaussian  $\varepsilon$ -graph model (56).

Based on the same methodology for parameter selection using K-fold cross validation (with K = 5) as described in Section 6.2.2, we compare the matrix recovery qualities of GRMC and the other two graph-agnostic matrix completion models. The results are shown in Table 1.

	SR=1%			SR=5%			SR=20%		
Method	MC MMMF GRMC			MC MMMF GRMC			MC MMMF GRMC		
GRALS1	0.0453	0.0351	0.0343	0.0778	0.0291	0.0272	0.0371	0.0246	0.0232
AltMin1	0.1218	0.0356	0.0344	0.0455	0.0332	0.0280	0.0317	0.0249	0.0244
AltMin2	0.1217	0.0352	0.0343	0.0455	0.0308	0.0275	0.0317	0.0247	0.0238
Euclidean GD	0.0455	0.0352	0.0343	0.0399	0.0299	0.0276	0.0261	0.0253	0.0241
Euclidean CG	0.0472	0.0350	0.0343	0.1443	0.0289	0.0272	0.0360	0.0246	0.0232
Qprecon RGD	0.0553	0.0351	0.0344	0.0477	0.0293	0.0273	0.0297	0.0246	0.0232
Qprecon RCG	0.0514	0.0350	0.0343	0.1875	0.0291	0.0271	0.0570	0.0246	0.0232
Qrightinv RGD	0.0454	0.0452	0.0349	0.0329	0.0326	0.0326	0.0300	0.0299	0.0300
Qrightinv RCG	0.0454	0.0452	0.0343	0.0812	0.0295	0.0273	0.0261	0.0254	0.0241

Table 1: Recovery scores (test RMSE) of the three matrix completion models under various sampling rates (SR): the unregularized matrix completion (MC), maximum-margin matrix factorization (MMMF) and graph-regularized matrix completion (GRMC).

Figure 8 shows the time efficiency of the methods tested in terms of the RMSE score per iteration.

#### 6.3.3 MovieLens dataset

The MovieLens 100K<sup>9</sup> dataset [19] consists of 10000 ratings (1 to 5) from 943 users on 1682 movies. Each user has rated at least 20 movies. The data was collected through the MovieLens web site (movielens.umn.edu) during the seven-month period from September 19th, 1997 through April 22nd, 1998. This data has been cleaned up—users who had less than 20 ratings or did not have complete demographic information were removed from this data set. For the graph-regularized matrix completion model, we construct a sparse row-wise similarity graph (on the set of users) with the Gaussian  $\varepsilon$ -graph model (56).

Based on the same methodology for parameter selection using K-fold cross validation (with K = 5) as described in Section 6.2.2, we compare the matrix recovery qualities of

 $<sup>^{8}</sup>$  https://archive.ics.uci.edu/ml/datasets/PEMS-SF

<sup>&</sup>lt;sup>9</sup>https://grouplens.org/datasets/movielens/100k/



Figure 8: Test RMSE per iteration on the *Traffic* dataset (m = 963, n = 10560). The rank chosen is for the model (9) is k = 18. (a): the sampling rate  $|\Omega|/mn = 1.0\%$ , (b):  $|\Omega|/mn = 5.0\%$ , (c):  $|\Omega|/mn = 20.0\%$ . In particular, the label (1s) of the dashed line refers to the method using backtracking-Armijo line search (Algorithm A.4).

GRMC and the other two graph-agnostic matrix completion models. The results are shown in Table 2. The RMSE scores of the GRMC model, returned by the methods tested using the selected parameter setting, are around 0.957, which is close to the RMSE score of 0.945 given by the graph-regularized method in [42] and is better than the scores of all other methods reported in [42]. Note that (i) the rank value chosen in the present experiment is the same as that in [42] and (ii) the graph Laplacian matrix used by Rao et al. [42] comes from side information, while the graph Laplacian matrix in the present experiment is constructed with the sparse  $\varepsilon$ -graph model (56), and (iii) in our experiment, the training set is 80% of the data entries in the ML100k dataset, while Rao et al. [42] used 90% of the available data. To achieve even better recovery scores under the GRMC framework, one needs to refine the construction of the graph Laplacian matrix either with models that are more adapted to the features of the data matrix or using more sensible user/movie-related information.

Methods	MC	MMMF	GRMC
GRALS1	2.076	0.984	0.957
GRALS2	1.203	0.983	0.957
Euclidean CG	1.411	0.986	0.957
AltMin1	4.069	0.984	0.956
AltMin2	4.018	0.984	0.956
Qprecon RGD	1.083	0.986	0.957
Qprecon RCG	1.917	0.984	0.959

Table 2: Matrix completion score (RMSE on test entries) of solutions to the three types of problem models: unregularized matrix completion (MC), Maximum-margin matrix factorization and Graph-regularized matrix completion (GRMC).

We also compare the time efficiency of the methods tested in terms of the RMSE score per iteration. Results are shown in Figure 9.



Figure 9: RMSE per iteration on the *MovieLens100k* dataset (m = 943, n = 1682). Rank parameter k = 10. The number of revealed entries is 80% of the 100k available ratings and the effective sampling rate  $|\Omega|/mn \approx 5.05\%$ . In particular, the label (1s) of the green line refers to the method using backtracking-Armijo line search (Algorithm A.4).

### 6.3.4 Discussion of real-data experiments

From both Table 1 and Table 2, we observe that the matrix recovery quality of solutions to the GRMC model (9) is superior to those of the other two graph-agnostic matrix completion models. From Figure 8 and Figure 9, we have the following observations:

- Our algorithms (Qprecon RGD, RCG) are faster than Euclidean GD and the baseline alternating minimization methods AltMin1, AltMin2.
- The time efficiency of Qprecon RCG and the state-of-the-art method GRALS1 are similar on the two real datasets tested. Observe that both Qprecon RCG and GRALS1 are considerably faster than the two AltMin methods, though GRALS1 and AltMin are based on the same alternating minimization strategy. This can be due to the programming language (C++ for GRALS1 and MATLAB for AltMin) and to GRALS's above-mentioned additional stopping criterion for the subproblem solver.
- The stepsize by line minimization (23) yields faster convergence behavior than backtracking line search (with respect to the Armijo rule, starting from an arbitrary guess  $s_0 = 1$  for the initial stepsize).

# 7 Conclusion

In this paper, we focused on a graph-regularized matrix factorization problem for matrix completion. We proposed efficient algorithms for the underlying optimization problem on the product space  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . Our proposed gradient descent and conjugate gradient methods are based on specially designed Riemannian metrics on  $\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  that are inspired from metrics on the Riemannian quotient manifold of fixed-rank matrices. Moreover, we focused on a stepsize selection method by exact line minimization, which results in a superior time efficiency compared to the approach using back-tracking line search. We provided rigorous theoretical analysis of the convergence property of the proposed Riemannian gradient descent algorithm.

We have investigated the matrix recovery qualities of various matrix completion models under various sampling rates: we found that the graph-based regularization does provide improvement for the matrix recovery quality compared to graph-agnostic matrix completion models, especially for relatively low sampling rates.

We have also conducted extensive experiments on synthetic data: we observed that our approach achieves significant speedup compared to several baseline methods, including a state-of-the-art method (GRALS) using alternating minimization, on various experimental settings. Moreover, we have shown via several tests that the proposed algorithms are much less influenced by changes in the initialization point or the scale of the data matrix. In our experiments on real-world data, we found that our methods produce solutions to the graphregularized matrix completion model in comparable or less time than the baseline and the state-of-the-art methods.

# Appendix A Algorithms

### A.1 Computation details of Algorithms 4.1–4.2 with line minimization

Algorithm A.I Computation of the International gradie	Algorithm	<b>4.1</b> Co	omputation	of the	Riemannian	gradient
---	-----------	---------------	------------	--------	------------	----------

**Input:**  $x = (G, H) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}, P_{\Omega}(M^{\star}) \in \mathbb{R}^{m \times n}, \Omega, \Theta^{\mathrm{r}} \in \mathbb{R}^{m \times m}, \Theta^{\mathrm{c}} \in \mathbb{R}^{n \times n}$ , and the parameter  $\alpha$ .

**Output:** Riemannian gradient  $\xi = (\xi_G, \xi_H) \in T_x (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}).$ 

1: Compute the residual  $S = P_{\Omega} (GH^T - M^{\star})$ .  $\# (2k+1)|\Omega|$  flops

2: Compute

 $\partial_G f(x) = SH + \alpha \Theta^{\mathsf{r}} G, \quad \partial_H f(x) = S^T G + \alpha \Theta^{\mathsf{c}} H.$ 

 $# 4(|\Omega| + \operatorname{nnz}(\Theta^{\mathrm{r}}) + \operatorname{nnz}(\Theta^{\mathrm{c}}))k$  flops

3: Compute

$$\xi = \left(\partial_G f(x) (G^T G + \delta I_k), \partial_H f(x) (H^T H + \delta I_k)\right) \quad w.r.t. (16)$$

 $# 4(m+n)k^2$  flops

or

$$\xi = \left(\partial_G f(x) (H^T H + \delta I_k)^{-1}, \partial_H f(x) (G^T G + \delta I_k)^{-1}\right) \quad w.r.t. \ (17)$$
  
# 4(m+n)k<sup>2</sup> + 2C<sub>chol</sub>k<sup>3</sup> flops, see (57)

**Computing the Riemannian gradient.** Detailed steps and their respective computational costs for computing the Riemannian gradient are given in Algorithm A.1. In the case of computing QPRECON (17): For the matrix inversion-related computations in the form of  $AB^{-1}$ , with  $A := \partial_G f(x) \in \mathbb{R}^{m \times k}$  and  $B := (G^T G + \delta I_k) \in \mathbb{R}^{k \times k}$ , a typical approach is to first take (once) a Cholesky decomposition of B, whose cost is  $C_{\text{chol}}k^3$ , and then compute the forward-and-backward substitution to get each of the m rows of  $AB^{-1}$ , which costs  $2mk^2$ . In brief, the flop counts of this line consists of

- Computing  $(G^T G + \delta I_k)$  and  $(H^T H + \delta I_k)$ :  $2(m+n)k^2$  flops,
- Computing the Cholesky decomposition of  $(G^T G + \delta I_k)$  and  $(H^T H + \delta I_k) : 2C_{\text{chol}}k^3$ , where  $C_{\text{chol}} = 1/3$ .
- Forward-and-backward substitutions:  $2(m+n)k^2$ ,

which sum to

$$4(m+n)k^2 + 2C_{\rm chol}k^3.$$
 (57)

The dominant term in (57) is  $4(m+n)k^2$  when  $k \ll \min(m, n)$ , which is the case for low-rank matrix approximation problems with a small rank parameter k and large data matrices.

The total number of flops needed for Algorithm A.1 is either of the following

$$(6k+1)|\Omega| + 4\mathrm{nnz}\,(\Theta)\,k + 4(m+n)k^2,\tag{58a}$$

$$(6k+1)|\Omega| + 4\mathrm{nnz}\,(\Theta)\,k + 4(m+n)k^2 + 2C_{\mathrm{chol}}k^3,\tag{58b}$$

where (58a) is for computing QRIGHTINV (16) and (58b) is for computing QPRECON (17). The dominant cost in Algorithm A.1 is for the computations in lines 1 and 2, which is

$$\mathcal{O}\left(\left(|\Omega| + \operatorname{nnz}\left(\Theta\right)\right)k\right),\$$

where we use the term  $\operatorname{nnz}(\Theta) := \operatorname{nnz}(\Theta^{\mathrm{r}}) + \operatorname{nnz}(\Theta^{\mathrm{c}})$  to denote the sum on the right-hand side, for simplicity. Indeed, when  $k \ll \min(m, n)$  and the sampling rate  $\rho$  is of the order of 10%, the terms  $(m+n)k \leq (m+n)k^2 \ll |\Omega|k = \rho mnk$ .

**Computing the cost function.** For the algorithms with the line search procedure (Algorithm A.4), the evaluation of the cost function is needed.

Algorithm A.2 Computation of the cost function (9)	
<b>Input:</b> $x = (G, H) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}, P_{\Omega}(M^{\star}) \in \mathbb{R}^{m \times n}, \Omega, \Theta^{\mathrm{r}} \in \mathbb{R}^{m \times n}$	$\mathbb{R}^{m \times m}, \Theta^{c} \in \mathbb{R}^{n \times n}$ , and the
parameter $\alpha$ .	
<b>Output:</b> Function value $f(x)$ in (9).	
1: Compute the residual $S = P_{\Omega} (GH^T - M^*)$ .	$\# (2k+1) \Omega $ flops
2: Compute $f_{\Omega}(x) := \frac{1}{2} \ S\ _{F}^{2}$ ,	$\# 2 \Omega $ flops
2. Compute $\operatorname{Pog}(n) := \operatorname{Tr}(C^T \operatorname{Or} C) + \operatorname{Tr}(H^T \operatorname{Oc} H)$	

3: Compute 
$$\operatorname{Reg}(x) := \operatorname{Tr} \left( G^{T} \Theta^{r} G \right) + \operatorname{Tr} \left( H^{T} \Theta^{c} H \right),$$
  
4: Return  $f(x) = f_{\Omega}(x) + \frac{\alpha}{2} \operatorname{Reg}(x).$   
 $\# \operatorname{2nnz} (\Theta) k + 2(m+n)k$  flops

Hence, the cost for evaluating once the objective function (9) is

$$FLOPS_{fobj} = (2k+3)|\Omega| + 2nnz(\Theta)k + 2(m+n)k.$$

$$(59)$$

To see the order of magnitude of the total cost: the dominant costs of Algorithm A.2 are  $\mathcal{O}((|\Omega| + \operatorname{nnz}(\Theta))k)$ . The total cost of Algorithm A.2 is  $\mathcal{O}((|\Omega| + \operatorname{nnz}(\Theta))k)$ .

Computing the conjugate gradient direction. The following schemes for computing the CG step parameter  $\beta_t$  in the Riemannian optimization setting (20) are adapted from nonlinear conjugate gradient schemes in the classical Euclidean setting, such as

Polak-Ribiere [41] (PR) 
$$\beta = \max\left(0, \frac{g_{x^t}\left(\xi^t - \xi^{t-1}, \xi^t\right)}{g_{x^t}\left(\xi^{t-1}, \xi^{t-1}\right)}\right),$$
 (60a)

Hestenes-Stiefel [22] (HS+) 
$$\beta = \max\left(0, \frac{g_{x^t}\left(\xi^t - \xi^{t-1}, \xi^t\right)}{g_{x^t}\left(\xi^t - \xi^{t-1}, \eta^{t-1}\right)}\right),$$
 (60b)

Fletcher-Reeves [15] (FR) 
$$\beta = \frac{g_{x^t}(\xi^t, \xi^t)}{g_{x^t}(\xi^{t-1}, \xi^{t-1})}.$$
 (60c)

A survey on nonlinear conjugate gradient can be found in [17]. Implementation of these schemes (60) can be found in the Riemannian optimization toolbox MANOPT [4]. In our experiments, we choose the modified Hestenes-Stiefel (HS+) rule. The flop counts for the HS+ rule is 5(m+n)k.

From Algorithm A.3, the total flop counts for computing once the Riemannian CG direction, given two consecutive Riemannian gradients, is

$$13(m+n)k. (61)$$

Algorithm A.3 Computation of the conjugate gradient direction

**Input:** iterates  $x^{t-1}, x^t \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ , gradients  $\xi^t \in T_x(\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ ,  $\xi^{t-1} \in T_{x^{t-1}}(\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ , previous CG direction  $\eta^{t-1} \in T_{x^{t-1}}(\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$ .

- **Output:** CG direction  $\eta^t \in T_x \left( \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k} \right).$
- 1: Compute: CG step parameter  $\beta_t$  with one of the schemes in (60) and then  $\eta^t = -\xi^t + \beta_t \eta^{t-1}$ . # 7(m+n)k flops
- 2: Compute the angle between the CG direction and the gradient:  $\theta = \langle \eta^t, \xi^t \rangle / \|\eta^t\| \|\xi^t\|.$ # 6(m+n)k flops
- 3: Reset to gradient if desired:  $\eta^t = \xi^t$  if  $\theta < 0.1$ .

**Computational cost of the line minimization** (23). This corresponds to the computations for  $c_1, ..., c_4$  in (66a)–(66d), which sums to  $(6k + 11)|\Omega| + 2nnz(\Theta)k + 4(m + n)k$ . Details are in Appendix A.4.

### A.2 Computational cost of Qprecon/Qrightinv RGD (linemin)

Each iteration of RGD (linemin) consists of (i) computing the stepsize by line minimization (23), the cost of which is in (67), (ii) conducting the descent step, the cost of which is (m + n)k flops, (iii) computing the new gradient, the cost of which is in (58), and (iv) computing the norm of the gradient, the cost of which is 2(m + n)k flops. Note that since the step size  $s_t$  is obtained by (23), which guarantees a sufficient decrease, there is no need for any additional line search steps. Therefore, the total flop counts for one iteration of Algorithm 4.1 is

$$12(k+1)|\Omega| + 6\mathrm{nnz}\,(\Theta)\,k + (m+n)(4k^2 + 7k).$$
(62)

### A.3 Computational cost of Qprecon/Qrightinv RCG (linemin)

RCG (linemin) needs to compute the nonlinear CG direction via Algorithm A.3, and its flop counts is larger than that of RGD (linemin) (62) by exactly that in (61). The total cost is

$$12(k+1)|\Omega| + 6\operatorname{nnz}(\Theta)k + 4(m+n)(k^2 + 5k).$$
(63)

### A.4 Stepsize computation via line minimization

Computing the stepsize (23) requires minimizing

$$f(G + s\eta_G, H + s\eta_H) - f(G, H)$$

for  $s \ge 0$ .

We have  $f(G + s\eta_G, H + s\eta_H) - f(G, H) = A + B$ , where

$$A = \frac{1}{2} \| P_{\Omega} \left( s(G\eta_H^T + \eta_G H^T) + s^2 \eta_G \eta_H^T \right) \|_F^2 + \left\langle P_{\Omega}(GH^T - M), P_{\Omega}(s(G\eta_H^T + \eta_G H^T) + s^2 \eta_G \eta_H^T) \right\rangle$$
(64)

and

$$B = \frac{1}{2} \operatorname{Tr} \left[ \left( sG^T L_r \eta_G + s\eta_G^T L_r G + s^2 \eta_G^T L_r \eta_G \right) + \left( sH^T L_c \eta_H + s\eta_H^T L_c H + s^2 \eta_H^T L_c \eta_H \right) \right].$$
(65)

These two equations lead to the following quartic polynomial form  $A + B = \sum_{j=1}^{4} c_j s^j$ , where

$$c_{1} = \left\langle P_{\Omega}(GH^{T} - M), P_{\Omega}(G\eta_{H}^{T} + \eta_{G}H^{T}) \right\rangle + \operatorname{Tr}\left(\eta_{G}^{T}L_{r}G + \eta_{H}^{T}L_{c}H\right), \quad (66a)$$

$$c_{2} = \frac{1}{2} \left\| P_{\Omega}(G\eta_{H}^{T} + \eta_{G}H^{T}) \right\|_{F}^{2} + \left\langle P_{\Omega}(GH^{T} - M), P_{\Omega}(\eta_{G}\eta_{H}^{T}) \right\rangle +$$

$$\frac{1}{2} \operatorname{Tr} \left( \eta_G^T L_r \eta_G + \eta_H^T L_c \eta_H \right), \tag{66b}$$

$$c_3 = P_{\Omega}(G\eta_H^T + \eta_G H^T), P_{\Omega}(\eta_G \eta_H^T),$$
(66c)

$$c_4 = \frac{1}{2} \| P_{\Omega}(\eta_G \eta_H^T) \|_F^2.$$
(66d)

The solution to  $s^*$  is selected from the real positive roots of the derivative of this quartic function, which is the polynomial of degree 3,  $(A+B)'(s) = \sum_{j=1}^{4} c_j s^{j-1}$ , the roots of which are easily computed.

**Computational costs.** In Algorithms 4.1–4.2, whenever the line minimization (23) is required, it always follows the computation of a Riemannian gradient, during which we have stored the following intermediate matrices (i)  $S = P_{\Omega}(GH^T - M^*) \in \mathbb{R}^{m \times n}$  and (ii)  $\Theta^{\mathrm{r}}G \in \mathbb{R}^{m \times k}, \Theta^{\mathrm{c}}H \in \mathbb{R}^{n \times k}$ . Hence, in the following list of flop counts, the computations related to the items above need not be counted:

- For  $c_1$  in (66a):  $(4k+3)|\Omega| + 2(m+n)k$  flops. Information stored:<sup>10</sup>  $P_{\Omega}(G\eta_H^T)$  and  $P_{\Omega}(\eta_G H^T)$ .
- For  $c_2$  in (66b):  $(2k+4)|\Omega| + 2\operatorname{nnz}(\Theta)k + 2(m+n)k$  flops. Information stored:  $P_{\Omega}(\eta_G \eta_H^T)$ .
- For  $c_3$  in (66c):  $2|\Omega|$  flops.
- For  $c_4$  in (66d):  $2|\Omega|$  flops.

These sum up to

$$(6k+11)|\Omega| + 2\mathrm{nnz}\,(\Theta)\,k + 4(m+n)k.$$
(67)

### A.5 The constant parameter $\delta$ in the definition of gradients

In all the experiments in Section 6, the gradients defined in (16) or (17) are used with a parameter  $\delta = 0$ . In this setting, the underlying metric (12) is not guaranteed to be positive definite and the metric (17) is not always well-defined at any iterate  $x \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$ . The convergence analysis does not cover the case where  $\delta = 0$  in (16)–(17). Nevertheless, we note that the convergence behavior of our proposed algorithms so far tested agrees with the theoretical results presented in Section 5. In fact, in all our experimental settings, the problem parameters ( $\alpha, \gamma_r, \gamma_c$ ) and the largest rank value k are chosen properly, especially that the rank parameter k is set to an underestimated value. Therefore, we did not observe any singularity in all the results presented in Section 6.

Figure 10 shows that the iterative results of the proposed RGD algorithms using a strictly positive  $\delta$  (for  $\delta$  set to  $10^{-4}$ ) and those using  $\delta = 0$  are almost the same. The experimental setting for this illustration is the same as in Section 6.2.2.

<sup>&</sup>lt;sup>10</sup>The information is stored only inside the current iteration.



Figure 10: RMSE per iteration on sythetic data. The data matrix  $M^*$  is generated via (48) and (53): Matrix size m = 1000, n = 900, rank  $r^* = 12$ . The rank parameter k = 8. The sampling rate  $|\Omega|/mn = 5\%$ . The results returned by the RGD algorithm using gradients defined by the two metrics with  $\delta = 10^{-4}$  (labeled with "-d") and  $\delta = 0$  respectively. The convergence behaviors of the algorithms with  $\delta = 0$  and with a small  $\delta > 0$  are almost same.

### A.6 Computation details of Algorithms 4.1–4.2 with Armijo line-search

The line search procedure by backtracking with respect to the Armijo rule (22) is given in Algorithm A.4.

Algorithm A.4 Armijo line search

**Input:**  $f : \mathcal{M} \mapsto \mathbb{R}$ , a descent direction a retraction  $\mathcal{R}$  on  $\mathcal{M}$ ,  $x^t \in \mathcal{M}$ , initial stepsize  $s_t^0 > 0$ and  $\sigma$ ,  $\beta \in ]0, 1[$ .

Output: s.

- 1: Initialize:  $s = s_t^0$ .
- 2: while  $f(x^t) f(\mathcal{R}_{x^t}(s\eta^t)) < \sigma s \langle -\operatorname{grad} f(x^t), \eta^t \rangle$  do
- 3:  $s \leftarrow \beta s$ .
- 4: end while

Computational cost of RGD (Armijo). RGD (Armijo) corresponds to Algorithm 4.1 using the backtracking line search with the Armijo condition at each iteration. Computing once the function value and the Riemannian gradient at the same time costs in total

$$FLOPS_{fobj} + FLOPS_{gradf} - [FLOPS(P_{\Omega}(GH^{T})) + FLOPS(\Theta^{r}G, \Theta^{c}H)]$$
  
=  $(6k+4)|\Omega| + 4nnz(\Theta)k + 4(m+n)k^{2} + 2(m+n)k.$  (68)

Hence, the computational cost of the t-th iteration is

$$(6k+4)|\Omega| + 4\mathrm{nnz}\,(\Theta)\,k + 4(m+n)k^2 + 2(m+n)k + n_t^{\mathrm{LS}}\mathrm{FLOPS}_{\mathrm{fobj}},\tag{69}$$

where  $FLOPS_{fobj}$  is defined in (59).

Flop counts for RCG (Armijo). RCG lsArmijo (Algorithm 4.2, with stepsizes chosen via the Armijo line search) needs to compute the nonlinear CG direction via Algorithm A.3, and its per-iteration cost is larger than that of RGD (Armijo) by the amount of (61). In sum, it is

$$(6k+4)|\Omega| + 4\mathrm{nnz}\,(\Theta)\,k + 12(m+n)k^2 + 17(m+n)k + n_t^{\mathrm{LS}}\mathrm{FLOPS}_{\mathrm{fobj}}.$$
 (70)

for an iteration  $t \ge 0$ .

### A.7 Algorithms using the Euclidean gradient

The cost for computing the Euclidean gradient is smaller than the cost of computing the variable metric gradient by the cost of line 3 of Algorithm A.1. Therefore, the computational cost per-iteration of Euclidean GD (linemin) is smaller than (62) by exactly  $4(m+n)k^2$ , which is

$$12(k+1)|\Omega| + 6nnz(\Theta)k + 7(m+n)k.$$
(71)

Computing the Euclidean CG step, using the same rule for computing the CG directions, requires the same cost as by (61), hence it equals

$$13(m+n)k. (72)$$

As a consequence, the computational cost per-iteration of Euclidean CG (linemin) is larger than (71) by exactly that of (72), which is

$$12(k+1)|\Omega| + 6\mathrm{nnz}\,(\Theta)\,k + 20(m+n)k.$$
(73)

### A.8 Computation details in GRALS [42]

This subsection contains a description of the algorithm proposed by Rao et al. [42] and a detailed list of flop counts for the standard CG steps involved in this algorithm. GRALS consists of the following two alternating least squares procedures

$$G^{t} = \arg\min_{G} f(G, H^{t-1}), \tag{74a}$$

$$H^t = \arg\min_{H} f(G^t, H), \tag{74b}$$

for a given initial point.

The quadratic forms of the least-squares systems (74a)–(74b) have the following structures. The subproblem (74a) is a least-squares problem whose objective  $f(G) := f(G, H^t)$ can be rewritten as a quadratic form of the vectorization of  $G^T$  via the identification  $f(G) = \tilde{f}_1(\operatorname{vec}(G^T)),$ 

$$\min_{s} \tilde{f}_{1}(s) = \frac{1}{2} s^{T} A^{(1)} s - \operatorname{vec}(H^{tT} P_{\Omega}(M^{\star})^{T})^{T} s,$$
(75)

where  $A^{(1)} \in \mathbb{R}^{km \times km}$  has the following structure,

$$A^{(1)} = \bar{B}^{(1)} + \alpha \Theta^{\mathrm{r}} \otimes I_k, \tag{76}$$

where  $\bar{B}^{(1)} \in \mathbb{R}^{km \times km}$  is block diagonal with *m* diagonal blocks  $(B_i^{(1)})_{i=1,\dots,m}$  of size  $k \times k$  such that

$$B_i^{(1)} = \sum_{j \in \Omega_i} h_j h_j^T, \tag{77}$$

for  $i \in [m]$ . Here the index sets  $\Omega_i := \{j \in [n] : (i, j) \in \Omega\}$  and

$$h_j = [H_{j1}, .., H_{jk}]^T \in \mathbb{R}^k$$
(78)

is the transpose of the j-th row of H.

Similarly, the subproblem (74b) can be solved by the same routines (Algorithm A.5 and A.6) as for (74a) by swapping the roles of G and H (and matrices in the regularization terms) in all computations of matrices involved. In the implementation of GRALS, the linear CG routine is used to solve the two subproblems in the form of (75). Algorithm A.6 is a standard CG descent procedure with  $A \in \mathbb{R}^{q \times q}$ ,  $b \in \mathbb{R}^{q}$  as inputs and  $x^{0}$  as the initial point. Note that GRALS [42] uses a warm-start scheme:  $x^{0}$  corresponds to latest iterate  $G^{t-1}$  (resp.  $H^{t-1}$ ) for the t-th step (74a) (resp. (74b)).

The Hessian-vector multiplication in the linear CG iteration (Algorithm A.6, line 13) is computed via Algorithm A.5.

# Algorithm A.5 Hessian-vector multiplication $A^{(1)}s$ [42]

**Input:** Data (known on  $\Omega$ )  $P_{\Omega}(M^{\star}) \in \mathbb{R}^{m \times n}, \Omega \subset [\![m]\!] \times [\![n]\!]$ . Quadratic form  $A^{(1)}$  in (76). Vector  $s := \operatorname{vec}(G^T) \in \mathbb{R}^{k \times m}$ . Laplacian-based matrix  $\overline{\Theta}$ . **Output:**  $A^{(1)}s$ . 1: **for** i = 1, ..., m **do** 2: Get  $g_i := [G_{i1}..., G_{ik}]^T$  from s (vectorization of  $G^T$ ). 3: Compute  $\tilde{g}_i = \sum_{j \in \Omega_i} h_j(h_j^T g_i)$ . # See (77). 4: **end for** 5: Get G from the vectorization  $s = \operatorname{vec}(G^T)$  and compute  $\tilde{G} = \overline{\Theta}G$ . 6: Return:  $\operatorname{vec}([\tilde{g}_1, ..., \tilde{g}_m]) + \operatorname{vec}(\tilde{G}^T)$ .

Algorithm A.6 CG Algorithm for solving (74b) (resp. (74a))

**Input:**  $A \in \mathbb{R}^{q \times q}$ , for q = nk (resp. mk), initial point  $x^0 \in \mathbb{R}^q$ . Accuracy parameter  $\epsilon$ , iteration budget  $n_{\rm CG}$ . **Output:**  $x^{\star} \in \mathbb{R}^{q}, n_{CG}^{\star}$ 1: Compute:  $b = \operatorname{vec}(P_{\Omega}(M^{\star})^T G) \in \mathbb{R}^q$  (resp.  $b = \operatorname{vec}(P_{\Omega}(M^{\star})H)$ ).  $\# 2|\Omega|k$  flops 2:  $r_0 = b - Ax^0$ . 3: for  $k = 0, ..., n_{CG}$  do Compute:  $||r_k||$ . # 2nk (resp. 2mk) flops 4: if  $||r_k|| \leq \epsilon ||r_0||$  then 5:Break; 6: end if 7: if k = 0 then 8: 9: $p_1 = r_0.$ else10: $p_{k+1} = r_k + \frac{\|r_k\|^2}{\|r_{k-1}\|^2} p_k.$ # 2nk (resp. 2mk) flops 11: end if 12:Compute:  $v_{k+1} = Ap_{k+1}$ . Compute:  $\beta = \frac{\|r_k\|^2}{p_{k+1}^T v_{k+1}}$ .  $\# 2(|\Omega| + \operatorname{nnz}(\Theta))k$  flops 13:# 2nk (resp. 2mk) flops 14: Compute:  $x_{k+1} = x_k + \beta p_{k+1}, r_{k+1} = r_k - \beta v_{k+1}.$ #4nk (resp. 4mk) flops 15:16: **end for** 17: Return  $x^* = x^k, n_{CG}^* = k$ .

Computational cost of GRALS. The number of flops required by Algorithm A.6 is:

$$2(n_{\rm CG}^{\star}+1)|\Omega|k+2n_{\rm CG}^{\star}nnz(\Theta)k+10n_{\rm CG}^{\star}nk$$
  
( resp.  $2(n_{\rm CG}^{\star}+1)|\Omega|k+2n_{\rm CG}^{\star}nnz(\Theta)k+10n_{\rm CG}^{\star}mk).$ 

During the t-th iteration in GRALS, let  $n_t^H$  (respectively  $n_t^G$ ) denote the number of CG iterations (*i.e.*  $n_{CG}^{\star}$  returned by this algorithm) required by Algorithm A.6 for solving the subproblem (74b) (resp. (74a)) at iteration t. Then the number of flops required by GRALS to complete the t-th iteration, from  $(G^t, H^t)$  to  $(G^{t+1}, H^{t+1})$  is

$$2(n_t^G + n_t^H + 2)|\Omega|k + 2(n_t^G + n_t^H) \operatorname{nnz}(\Theta)k + 10(n_t^G m + n_t^H n)k.$$
(79)

Figure 11 shows the RMSE per iteration, where the x-axis is represented either by the wall time recorded at each iteration or the cumulative cost (in flops) required by the main computational steps in each of the algorithms at each iteration.



Figure 11: RMSE per iteration. Left: the x-axis is wall time at each iteration. Right: the x-axis is the cumulative computational cost at each iteration. Experimental settings:  $m = 1000, n = 900, k = r^* = 10, |\Omega|/mn = 20.0\%$ .  $M^*$  is generated with the model (50) and is partially observed without any noise.

### A.8.1 Our implementation (AltMin)

In addition to GRALS [42], we implement the alternating minimization method ((74a)–(74b)) with a linear CG solver that is controlled by the two following parameters,

- $n_{\rm CG}$ : the iteration budget for each of the two least-squares subproblems.
- $\epsilon$ : tolerance parameter to control the accuracy of the solutions to each of the two subproblems.

The most costly computation in AltMin/GRALS is the computation of  $A^{(1)}$ vec  $(G^T)$ and  $A^{(2)}$ vec  $(H^T)$ . Algorithm A.7 avoids searching for indices in the subset  $\Omega_i$  for each  $i \in [m]$  by using the following incremental procedure. The notations therein are adapted to the computation of  $B^{(1)}$ vec  $(G^T)$ . Note that for the computation of  $B^{(2)}$ vec  $(H^T)$ , this algorithm applies by swapping the roles of G and H.

### Algorithm A.7 Hessian-vector multiplication $B^{(1)}s$

Input:  $\Omega \subset \llbracket m \rrbracket \times \llbracket n \rrbracket$ .  $H \in \mathbb{R}^{n \times k}$ , vector  $s := \operatorname{vec}(G^T) \in \mathbb{R}^{km}$ . Output:  $B^{(1)}s$ , for  $B^{(1)}$  in (77). 1: Initialize the k-dimensional vectors:  $y_i = \mathbf{0}$  for i = 1, ..., m. 2: for  $l = 1, ..., |\Omega|$  do 3: Get  $(i_l, j_l)$ : the l-th pair of  $\Omega$ . 4: Get  $h_{j_l}$  from H. 5: Get  $g_{i_l} = [G_{i_l1}, ..., G_{i_lk}]^T$ , which is  $(s_{i_lk-k+1}, ..., s_{i_lk})$ . 6: Compute  $y_{i_l} = y_{i_l} + h_{j_l}(h_{j_l}^T g_{i_l})$ . 7: end for 8: Return:  $\operatorname{vec}([y_1, ..., y_m]) \in \mathbb{R}^{km}$ .

### **A.9** The Hessian of the objective function of (9)

The second-order Euclidean directional derivative of f at  $x = (G, H) \in \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k}$  along a direction  $\xi = (\xi_G, \xi_H) \in T_x (\mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k})$  is defined as

$$\nabla^2 f(x)[\xi] := \frac{d}{dt} \nabla f(x + t\xi)|_{t=0}.$$
(80)

The gradient vector field has the following expression,

$$\nabla f(x) = \left(SH + \alpha \Theta^{\mathrm{r}} G, S^{T} G + \alpha \Theta^{\mathrm{c}} H\right), \qquad (81)$$

where  $S := P_{\Omega}(GH^T - M)$ . To simplify notations, we calculate the two matrix components separately,

$$\frac{d}{dt}\partial_{G}f(x+t\xi)|_{t=0} = \lim_{t\to 0} \frac{1}{t} \Big[ P_{\Omega}((G+t\xi_{G})(H+t\xi_{H})^{T} - M)(H+t\xi_{H}) - P_{\Omega}(GH^{T} - M)H + t\alpha\Theta^{r}\xi_{G} \Big] 
= P_{\Omega}(G\xi_{H}^{T} + \xi_{G}H^{T})H + S\xi_{H} + \alpha\Theta^{r}\xi_{G}.$$
(82)
  
(82)

Similarly,  $\frac{d}{dt}\partial_H f(x+t\xi)|_{t=0} = P_{\Omega}(G\xi_H^T + \xi_G H^T)^T G + S^T \xi_G + \alpha \Theta^c \xi_H$ . Hence we have

$$\nabla^2 f(x)[\xi] = \begin{pmatrix} P_{\Omega}(G\xi_H^T + \xi_G H^T)H + S\xi_H + \alpha \Theta^{\mathrm{r}}\xi_G \\ P_{\Omega}(G\xi_H^T + \xi_G H^T)^T G + S^T \xi_G + \alpha \Theta^{\mathrm{c}}\xi_H \end{pmatrix}.$$
(84)

# Appendix B The matrix model with graph information and the graph Laplacian-based regularization

The model (48) is of particular interest because it models a wide range of real data matrices whose entries present pairwise similarities. Figure 1 shows differences between Z (48a) and  $A^{\rm r}Z$  (48b) in both the data entries and in the graph spectral domain. By using the concept of graph Fourier transforms (*e.g.* [46]), we illustrate how ( $g, A^{\rm r}, A^{\rm c}$ ) in (49) transforms a Gaussian random matrix Z into a matrix with more apparent pairwise similarities on the given graphs.

By definition (e.g. [46]), the graph Fourier transform of a vector  $f \in \mathbb{R}^m$  with respect to the graph Laplacian  $L^{\mathbf{r}} = U\Lambda U^T$ , is

$$\widehat{f}(\lambda_l) = (Ue_l)^T f, \forall l \in \llbracket m \rrbracket$$

Now we compare the smoothness of the Gaussian low-rank model  $Z = FQ^T$  in (48a) and the graph-based model  $X = A^r Z$  (48b) with respect to the row-wise similarity graph  $L^r$ : For any  $j = 1, \ldots, k$ , the graph Fourier coefficients of the *j*-th column of the Gaussian random matrix  $F \in \mathbb{R}^{m \times k}$  and transformed matrix  $G = A^r F$  are

$$\widehat{F_{(j)}}(\lambda_l) = (U^T e_l)^T F_{(j)},$$

$$\widehat{G_{(j)}}(\lambda_l) = (U^T e_l)^T A^{\mathrm{r}} F_{(j)} = \sqrt{g(\lambda_l)} e_l^T F_{(j)}, \forall l \in \llbracket m \rrbracket,$$

where  $F_{(j)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_F^2 I_m)$ . From basic calculations, the amplitudes of their graph Fourier coefficients satisfy

$$\mathbb{E}\left[|\widehat{F_{(j)}}(\lambda_l)|^2\right] = \sum_{i=1}^m U_{(i,l)}^2 \mathbb{E}\left[F_{(i,j)}^2\right] = \sigma_F^2 \sum_{i=1}^m U_{(i,l)}^2 = \sigma_F^2,$$
(85)

$$\mathbb{E}\left[|\widehat{G_{(j)}}(\lambda_l)|^2\right] = g(\lambda_l)\mathbb{E}\left[\|e_l^T F_{(j)}\|_2^2\right] = \sigma_F^2 g(\lambda_l)$$
(86)

Therefore, when g is decreasing on  $(0, +\infty)$  as defined in (50), the sequence  $(g(\lambda_l))_l$  is decreasing for increasing values of  $(\lambda_l)_{l=2,...,m}$ . Note that a small eigenvalue  $\lambda$  corresponds eigenfunctions on the graph with small variations. This means the energy of  $G_{(j)}$  in the graph Fourier domain, determined by  $(|\widehat{G_{(j)}}(\lambda_l)|)_{1 \le l \le m}$ , is mostly concentrated on the "low graph-vertex frequencies".

The overall variations of the matrix factor G is related to  $\text{Tr}(G^T L^r G)$  in the regularizer of our main problem (9) as follows,

$$\frac{1}{k} \operatorname{Tr} \left( G^T L^r G \right) = \frac{1}{k} \sum_{j=1}^k \sum_{l=1}^m \lambda_l^2 |\widehat{G_{(j)}}(\lambda_l)|^2 = \sum_{l=1}^m \lambda_l^2 \widetilde{\mathbb{E}} |\widehat{G_{(1)}}(\lambda_l)|^2.$$
(87)

The weighted-sum expression (87) dictates that  $\operatorname{Tr}(G^T L^r G)$  is small when the amplitudes  $(\|\widehat{G}_{(1)}(\lambda_l)\|)_l$  are concentrated on low-frequencies, such as in (86). The same property applies to the factor H with respect to  $L^c$ . This reflects that the graph-based regularizer

$$S_L(x) = \operatorname{Tr}\left(G^T L^{\mathrm{r}} G\right) + \operatorname{Tr}\left(H^T L^{\mathrm{c}} H\right),$$

quantifies the smoothness of the entries of (G, H) on the row and column index sets with respect to the row-wise and column-wise similarity graphs ( $\mathcal{G}^{r}$  and  $\mathcal{G}^{c}$ ), as explained in (2).

### Acknowledgement

Dedicated to Paul Van Dooren on the occasion of his 70th birthday. His many fundamental contributions to numerical computation have been an inspiration to us.

# References

- M. Belkin, I. Matveeva, and P. Niyogi. Tikhonov regularization and semi-supervised learning on large graphs. In 2004 IEEE International Conference on Acoustics, Speech, and Signal Processing, volume 3, pages iii–1000. IEEE, 2004.
- [2] M. Belkin and P. Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*, 15(6):1373–1396, 2003.

- [3] N. Boumal, P. A. Absil, and C. Cartis. Global rates of convergence for nonconvex optimization on manifolds. *IMA Journal of Numerical Analysis*, 39(1):1–33, 2019.
- [4] N. Boumal, B. Mishra, P.-A. Absil, and R. Sepulchre. Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15:1455–1459, 2014.
- [5] E. Candès and B. Recht. Exact low-rank matrix completion via convex optimization. 2008 46th Annual Allerton Conference on Communication, Control, and Computing, (m):1–49, 2008.
- [6] E. J. Candès and B. Recht. Exact Matrix Completion via Convex Optimization. Foundations of Computational Mathematics, 9(6):717–772, 2009.
- [7] E. J. Candès and T. Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- [8] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky. The Convex Geometry of Linear Inverse Problems. *Foundations of Computational Mathematics*, 12(6):805–849, 2012.
- B. Chazelle. An improved algorithm for the fixed-radius neighbor problem. Information Processing Letters, 16(4):193–198, 1983.
- [10] J. Chen, H. A. Gov, and Y. Saad. Fast Approximate kNN Graph Construction for High Dimensional Data via Recursive Lanczos Bisection Haw-ren Fang. *Journal of Machine Learning Research*, 10, 2009.
- [11] F. R. K. Chung. Spectral Graph Theory. American Mathematical Society, 1997.
- [12] R. R. Coifman, S. Lafon, A. B. Lee, M. Maggioni, B. Nadler, F. Warner, and S. W. Zucker. Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. *Proceedings of the national academy of sciences*, 102(21):7426–7431, 2005.
- [13] R. R. Coifman and M. Maggioni. Diffusion wavelets. Applied and Computational Harmonic Analysis, 21(1):53–94, 2006.
- [14] S. Dong, P.-A. Absil, and K. A. Gallivan. Preconditioned Conjugate Gradient Algorithms for Graph Regularized Matrix Completion. In European Symposium on Artificial Neural Networks (ESANN), pages 239–244, 2019.
- [15] R. Fletcher and C. M. Reeves. Function minimization by conjugate gradients. The Computer Journal, 7(2):149–154, 1964.
- [16] R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. 34th International Conference on Machine Learning, ICML 2017, 3:1990–2028, 2017.
- [17] W. W. Hager and H. Zhang. A Survey of Nonlinear Conjugate Gradient Methods. *Pacific journal of Optimization*, 2(1):35–58, 2006.
- [18] M. Hardt. Understanding alternating minimization for matrix completion. Proceedings -Annual IEEE Symposium on Foundations of Computer Science, FOCS, pages 651–660, 2014.

- [19] F. M. Harper and J. A. Konstan. The MovieLens Datasets: History and Context. ACM Transactions on Interactive Intelligent Systems (TiiS) 5, 4, Article 19 (December 2015), 19 pages., 2015.
- [20] T. Hastie, R. Mazumder, J. D. Lee, and R. Zadeh. Matrix completion and low-rank SVD via fast alternating least squares. *The Journal of Machine Learning Research*, 16(1):3367–3402, 2015.
- [21] X. He and P. Niyogi. Locality preserving projections. In Advances in neural information processing systems, pages 153–160, 2004.
- [22] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. Journal of research of the National Bureau of Standards, 49:409–436, 1952.
- [23] P. Jain and I. S. Dhillon. Provable Inductive Matrix Completion. Technical report, 2013.
- [24] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. Proceedings of the 45th annual ACM symposium on Symposium on theory of computing - STOC '13, page 665, 2013.
- [25] G. James, D. Witten, T. Hastie, and R. Tibshirani. An introduction to statistical learning, volume 112. Springer, 2013.
- [26] V. Kalofolias, X. Bresson, M. Bronstein, and P. Vandergheynst. Matrix Completion on Graphs. In NIPS2014 - Robustness in High Dimension, 2014.
- [27] R. H. Keshavan, A. Montanari, and S. Oh. Matrix completion from a few entries. *IEEE transactions on information theory*, 56(6):2980–2998, 2010.
- [28] R. H. Keshavan and S. Oh. OptSpace : A Gradient Descent Algorithm on the Grassman Manifold for Matrix Completion. 2009.
- [29] C. Ma, K. Wang, Y. Chi, and Y. Chen. Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In 35th International Conference on Machine Learning, ICML 2018, volume 8, pages 5264–5331, 2018.
- [30] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11(Jan):19–60, 2010.
- [31] R. Mazumder, T. Hastie, H. Edu, R. Tibshirani, T. Edu, and T. Jaakkola. Spectral Regularization Algorithms for Learning Large Incomplete Matrices. *Journal of Machine Learning Research*, 11:2287–2322, 2010.
- [32] G. Meyer, S. Bonnabel, and R. Sepulchre. Linear regression under fixed-rank constraints: a Riemannian approach. In *Proceedings of the 28th international conference on machine learning*, 2011.
- [33] B. Mishra, K. A. Apuroop, and R. Sepulchre. A Riemannian geometry for low-rank matrix completion. arXiv preprint arXiv:1211.1550, 2012.
- [34] B. Mishra, G. Meyer, F. Bach, and R. Sepulchre. Low-rank optimization with trace norm penalty. SIAM Journal on Optimization, Society for Industrial and Applied Mathematics, 23(4):2124–2149, 2013.

- [35] B. Mishra, G. Meyer, S. Bonnabel, and R. Sepulchre. Fixed-rank matrix factorizations and Riemannian low-rank optimization. *Computational Statistics*, 29(3-4):591– 621, 2014.
- [36] Y. Nesterov. Introductory Lectures on Convex Optimization, volume 87. Springer Publishing Company, Incorporated, 1 edition, 2004.
- [37] T. Ngo and Y. Saad. Scaled gradients on Grassmann manifolds for matrix completion. In Advances in Neural Information Processing Systems, pages 1412–1420, 2012.
- [38] L. T. Nguyen, J. Kim, and B. Shim. Low-Rank Matrix Completion: A Contemporary Survey. *IEEE Access*, 7:94215–94237, jul 2019.
- [39] B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? Vision Research, 37(23):3311–3325, 1997.
- [40] N. Perraudin, J. Paratte, D. Shuman, L. Martin, V. Kalofolias, P. Vandergheynst, and D. K. Hammond. GSPBOX: A toolbox for signal processing on graphs. ArXiv e-prints, aug 2014.
- [41] E. Polak and G. Ribiere. Note sur la convergence de méthodes de directions conjuguées. Rev. Francaise Informat Recherche Opertionelle, pages 35–43, 1969.
- [42] N. Rao, H.-F. Yu, P. Ravikumar, and I. S. Dhillon. Collaborative Filtering with Graph Information: Consistency and Scalable Methods. In Advances in Neural Information Processing Systems 28, pages 2107–2115. 2015.
- [43] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. SIAM Review, 52(3):471–501, 2010.
- [44] J. D. M. Rennie and N. Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *Proceedings of the 22nd international conference on Machine learning - ICML '05*, pages 713–719, New York, New York, USA, 2005. ACM Press.
- [45] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst. The emerging field of signal processing on graphs. *IEEE Signal Processing Magazine*, 30(3):83–98, 2013.
- [46] D. I. Shuman, B. Ricaud, and P. Vandergheynst. Vertex-frequency analysis on graphs. Applied and Computational Harmonic Analysis, 40(2):260–291, mar 2016.
- [47] D. Spielman. Spectral Graph Theory. Combinatorial Scientific Computing. Chapman and Hall/CRC Press, 2012.
- [48] N. Srebro, J. D. Rennie, and T. S. Jaakkola. Maximum-margin matrix factorization. Advances in Neural Information Processing Systems, 17:1329–1336, 2005.
- [49] R. Sun and Z. Q. Luo. Guaranteed Matrix Completion via Non-Convex Factorization. IEEE Transactions on Information Theory, 62(11):6535–6579, 2016.
- [50] A. Uschmajew and B. Vandereycken. Geometric methods on low-rank matrix and tensor manifolds. Variational methods for nonlinear geometric data and applications (P. Grohs, M. Holler, A. Weinmann, eds.). Springer, 2020.
- [51] B. Vandereycken. Low-Rank Matrix Completion by Riemannian Optimization. SIAM Journal on Optimization, 23(2):1214–1236, 2013.

- [52] Y.-X. Wang and Y.-J. Zhang. Nonnegative matrix factorization: A comprehensive review. IEEE Transactions on Knowledge and Data Engineering, 25(6):1336–1353, 2012.
- [53] M. Xu, R. Jin, and Z.-H. Zhou. Speedup matrix completion with side information: Application to multi-label learning. In Advances in neural information processing systems, pages 2301–2309, 2013.
- [54] Y. Xu and W. Yin. A Block Coordinate Descent Method for Regularized Multiconvex Optimization with Applications to Nonnegative Tensor Factorization and Completion. *SIAM Journal on Imaging Sciences*, 6(3):1758–1789, 2013.
- [55] H.-F. Yu, N. Rao, and I. S. Dhillon. Temporal Regularized Matrix Factorization for High-dimensional Time Series Prediction. In Advances in Neural Information Processing Systems 29, pages 847–855, 2016.
- [56] F. Zhang and E. R. Hancock. Graph spectral image smoothing using the heat kernel. *Pattern Recognition*, 41(11):3328–3342, 2008.
- [57] H. Zhang and S. Sra. First-order methods for geodesically convex optimization. In Conference on Learning Theory, pages 1617–1638, 2016.
- [58] X. Zhang, S. Du, and Q. Gu. Fast and Sample Efficient Inductive Matrix Completion via Multi-Phase Procrustes Flow. In J. Dy and A. Krause, editors, *Proceedings of* the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pages 5756–5765, Stockholmsmässan, Stockholm Sweden, 2018. PMLR.
- [59] Z. Zhao, L. Zhang, X. He, and W. Ng. Expert Finding for Question Answering via Graph Regularized Matrix Completion. *IEEE Transactions on Knowledge and Data Engineering*, 27(4):993–1004, apr 2015.
- [60] G. Zhou, W. Huang, K. A. Gallivan, P. Van Dooren, and P. A. Absil. A Riemannian rank-adaptive method for low-rank optimization. *Neurocomputing*, 192:72–80, jun 2016.
- [61] T. Zhou, H. Shan, A. Banerjee, and G. Sapiro. Kernelized probabilistic matrix factorization: Exploiting graphs and side information. In *Proceedings of the 2012 SIAM international Conference on Data mining*, pages 403–414. SIAM, 2012.