



HAL
open science

Spurious Valleys, Spurious Minima and NP-hardness of Sparse Matrix Factorization With Fixed Support

Quoc-Tung Le, Elisa Riccietti, Rémi Gribonval

► **To cite this version:**

Quoc-Tung Le, Elisa Riccietti, Rémi Gribonval. Spurious Valleys, Spurious Minima and NP-hardness of Sparse Matrix Factorization With Fixed Support. 2021. hal-03364668v3

HAL Id: hal-03364668

<https://hal.science/hal-03364668v3>

Preprint submitted on 29 Nov 2021 (v3), last revised 16 Nov 2022 (v5)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

SPURIOUS VALLEYS, SPURIOUS MINIMA AND NP-HARDNESS OF SPARSE MATRIX FACTORIZATION WITH FIXED SUPPORT

QUOC-TUNG LE*, ELISA RICCIETTI*, AND REMI GRIBONVAL*

Abstract. The problem of approximating a dense matrix by a product of sparse factors is a fundamental problem for many signal processing and machine learning tasks. It can be decomposed into two subproblems: finding the position of the non-zero coefficients in the sparse factors, and determining their values. While the first step is usually seen as the most challenging one due to its combinatorial nature, this paper focuses on the second step, referred to as sparse matrix approximation with fixed support. First, we show its NP-hardness, while also presenting a nontrivial family of supports making the problem practically tractable with a dedicated algorithm. Then, we investigate the landscape of its natural optimization formulation, proving the absence of spurious local valleys and spurious local minima, whose presence could prevent local optimization methods to achieve global optimality. The advantages of the proposed algorithm over state-of-the-art first-order optimization methods are discussed.

Key words. Sparse Matrix Factorization, Fixed Support, NP-hardness, Landscape

AMS subject classifications. 15A23, 90C26

1. Introduction. Matrix factorization with sparsity constraints is the problem of approximating a (possibly dense) matrix as the product of two or more sparse factors. This is playing an important role in many domains and applications such as dictionary learning and signal processing [21, 19, 18], linear operator acceleration [13, 12], deep learning [2], to mention only a few.

In this work, we consider a particular instance of the matrix factorization problem with sparsity constraints, in which just two factors are considered and they have prescribed supports. We call this problem *fixed support (sparse) matrix factorization* (FSMF). In details, given a matrix $A \in \mathbb{R}^{m \times n}$, we look for two sparse factors X, Y that solve the following problem:

$$\begin{aligned}
 \text{(FSMF)} \quad & \underset{X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}}{\text{Minimize}} && L(X, Y) = \|A - XY^\top\|^2 \\
 & \text{Subject to:} && \text{supp}(X) \subseteq I \text{ and } \text{supp}(Y) \subseteq J
 \end{aligned}$$

where $\|\cdot\|$ is the Frobenius norm, $I \subseteq \llbracket m \rrbracket \times \llbracket r \rrbracket$, $J \subseteq \llbracket n \rrbracket \times \llbracket r \rrbracket^1$ are given support constraints, i.e., $\text{supp}(X) \subseteq I$ implies that $\forall (i, j) \notin I, X_{ij} = 0$.

The main aim of this work is to investigate the theoretical properties of (FSMF). To the best of our knowledge the analysis of matrix factorization problems with fixed supports has never been addressed in the literature. This analysis is however interesting, for at least two reasons.

Firstly, there are many practical applications in which the solution of this problem is required. Indeed, there are matrices that can be written as the product of factors whose support is known in advance. This is the case for instance of many fast transforms such as the Discrete Fourier Transform (DFT) or the the Hadamard Transform (HT), in which the fixed supports of the factors have the butterfly structure [12, 2].

Moreover, (FSMF) can be seen as a subproblem of a more general matrix factor-

*Univ Lyon, ENS de Lyon, UCBL, CNRS, Inria, LIP, F-69342, LYON Cedex 07, France (quoc-tung.le@ens-lyon.fr, elisa.riccietti@ens-lyon.fr, remi.gribonval@inria.fr).

¹ $\forall m \in \mathbb{N}, \llbracket m \rrbracket := \{1, \dots, m\}$

ization problem with structured sparsity constraints:

$$(1.1) \quad \begin{aligned} & \underset{X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}}{\text{Minimize}} && L(X, Y) = \|A - XY^\top\|^2 \\ & \text{Subject to:} && X \in \Sigma_X \text{ and } Y \in \Sigma_Y \end{aligned}$$

where $\Sigma_X \subseteq \mathbb{R}^{m \times r}$, $\Sigma_Y \subseteq \mathbb{R}^{n \times r}$ are some sets of structured sparse matrices. Relevant examples of such sets are for instance the sets of matrices with:

- at most k non-zero coefficients $\Sigma_k^{\text{total}} = \{X \in \mathbb{R}^{m \times r} \mid \|X\|_0 \leq k\}$;
- at most k non-zero coefficients per column (resp. per row) $\Sigma_k^{\text{col}} = \{X \in \mathbb{R}^{m \times r} \mid \|X_{\bullet, i}\|_0 \leq k, \forall i = 1, \dots, r\}$ (resp. $\Sigma_k^{\text{row}} = \{X \in \mathbb{R}^{m \times r} \mid \|X_{i, \bullet}\|_0 \leq k, \forall i = 1, \dots, m\}$),

where for a vector or matrix X , $\|X\|_0$ counts the number of nonzero entries in X .

Any heuristic algorithm for the solution of (1.1) will eventually need to deal with a subproblem of the form (FSMF), one way or another. Indeed, matrix factorization with sparsity constraints somehow generalizes the sparse recovery problem [4], in which we want to recover a sparse vector $x \in \mathbb{R}^n$ from the knowledge of its measurement vector (possibly corrupted by noise) $y = Ax \in \mathbb{R}^m$ with known measurement matrix $A \in \mathbb{R}^{m \times n}$. Mimicking the decomposition of the classical sparse recovery problem into a support recovery step and a coefficient recovery step, Problem (1.1) can also be split into two subproblems:

- 1) Determine the supports of X and Y , i.e. the set of indices $\text{supp}(X)$, $\text{supp}(Y)$ whose coefficients are different from zero. For instance, if $\Sigma_X = \Sigma_Y = \Sigma_k^{\text{total}}$, we need to identify the position of (at most) k non-zero coefficients of X and Y .
- 2) Determine the value of the coefficients in the supports of X and Y .

The solution of a problem in the form of (FSMF) will be needed both for *one-step* algorithms that jointly estimate the supports and coefficients, and for the *two-step* algorithms that solve the two problems successively. Also, as it happens in sparse linear regression, many common post-processing methods consist in "debiasing" the solution by a two-step approach [2].

Our aim is to then study the theoretical properties of (FSMF) and in particular to assess its difficulty. Assessing the difficulty of this subproblem is crucial to have a good understanding also of the difficulty of the full problem (1.1).

In particular, we consider three complementary aspects related to (FSMF).

First, we show the NP-hardness of (FSMF). While this result contrasts with the theory established for coefficient recovery with a fixed support in the classical sparse recovery problem (that can be trivially addressed by least squares), it is in line with the known hardness of related matrix factorization with additional constraints or different losses. Indeed, famous variants of matrix factorization such as non-negative matrix factorization (NMF) [23, 20], weighted low rank [5] and matrix completion [5] were all proved to be NP-hard. We prove the NP-hardness by reduction from the Matrix Completion problem with noise. To our knowledge this proof is new and cannot be trivially deduced from any existing result on the more classical full support case (i.e., the case in which $I = \llbracket m \rrbracket \times \llbracket r \rrbracket$, $J = \llbracket n \rrbracket \times \llbracket r \rrbracket$, which is equivalent to low rank matrix approximation [3]).

Second, we show that despite the hardness of (FSMF) in the general case, many pairs of support constraints (I, J) make the problem solvable by an effective direct algorithm based on the block singular value decomposition (SVD). The investigation of those supports is also covered in this work and a dedicated polynomial algorithm is proposed to deal with this family of supports. This includes for example the full support case. Our analysis of tractable instances of (FSMF) actually includes and substantially

generalizes the analysis of the instances that can be classically handled with the SVD decomposition. In fact, the presence of the constraints on the support makes it impossible to directly use the SVD to solve the problem, because coefficients outside the support have to be zero. However, the presented family of support constraints allows for an iterative decomposition of the problem into "blocks" that can be exploited to build up an optimal solution using blockwise SVDs. This technique can also be seen in many sparse representations of matrices (for example, \mathcal{H} -matrices or hierarchical matrices [7, 8]) to allow fast matrix-vector and matrix-matrix multiplication. In fact, matrices admitting hierarchical and related structures as in [7, 8] can in many cases be written as the product of sparse matrices with known supports. Our consideration of (FSMF) can be viewed as a generalization of this approach.

The third contribution of this paper is the study of the landscape of function L of (FSMF). Notably, we investigate the existence of *spurious local minima* and *spurious local valleys*, which will be collectively referred to as *spurious objects*. They will be formally introduced in Section 4, but intuitively these objects may represent a challenge for the convergence of local optimization methods.

The landscape of the loss functions for neural networks in general, and for linear neural networks in particular, has been a popular subject of study recently. In particular, great attention has been devoted to the investigation of the properties of critical points (i.e., points where the gradient vanishes) and global optima of the training problem with quadratic loss [9, 25, 14]. These works have direct link to ours since matrix factorization (without any constraint) can be seen as a specific case of neural network (with two layers, no bias and linear activation function).

Notably it has been proved [25] that for linear neural networks, every local minimum is a global minimum and if the network is shallow (i.e., there is only one hidden layer), critical points are either global minima or strict saddle points (i.e., their Hessian have at least one *strictly*-negative eigenvalue). However, there is still a *tricky* type of landscape that could represent a challenge for local optimization methods and has not been covered until recently: spurious local valleys [15, 24].

To the best of our knowledge, existing analyses of spurious local valleys are proposed for matrix factorization problems without support constraints, cf. [25, 24, 9], while the study of the landscape of (FSMF) remains untouched in the literature and our work can be considered as a generalization of such previous results.

To summarize, our main contributions in this paper are:

- 1) We prove that (FSMF) is NP-hard in Theorem 2.4.
- 2) We introduce families of support constraints (I, J) making (FSMF) tractable (Theorem 3.3 and Theorem 3.8) and provide dedicated polynomial algorithms for those families.
- 3) We show that the landscape of (FSMF) corresponding to the support pairs (I, J) in these families are free of spurious local valleys, regardless of the factorized matrix A (Theorem 4.12, Theorem 4.13). We also investigate the presence of spurious local minima for such families (Theorem 4.12, Theorem 4.19).
- 4) These results might suggest a conjecture that holds true for the full support case: an instance of (FSMF) is tractable if and only if their corresponding landscape is benign, i.e. free of spurious objects. We give a counter-example to this conjecture (Example 4.22) and show experimentally that first-order methods for the fixed support matrix factorization problem can fail despite a benign landscape and that a good initialization is really important.

1.1. Notations. For $n \in \mathbb{N}$, define $\llbracket n \rrbracket := \{1, \dots, n\}$. The notation $\mathbf{0}$ (resp. $\mathbf{1}$) stands for a matrix with all zeros (resp. all ones) coefficients. The identity matrix of size $n \times n$ is denoted by \mathbf{I}_n . Given a matrix $A \in \mathbb{R}^{m \times n}$ and $T \subseteq \llbracket n \rrbracket$, $A_{\bullet, T} \in \mathbb{R}^{m \times |T|}$ is the submatrix of A restrained to the columns indexed in T while $A_T \in \mathbb{R}^{m \times n}$ is the matrix that has the same columns as A for indexes in T and is zero elsewhere. If $T = \{k\}$ is a singleton, $A_{\bullet, T}$ is simplified as $A_{\bullet, k}$ (the k^{th} column of A). For $(i, j) \in \llbracket m \rrbracket \times \llbracket n \rrbracket$, $A_{i, j}$ is the coefficient of A at index (i, j) . If $S \subseteq \llbracket m \rrbracket$, $T \subseteq \llbracket n \rrbracket$, then $A_{S, T} \in \mathbb{R}^{|S| \times |T|}$ is the submatrix of A restrained to rows and columns indexed in S and T respectively.

A support constraint I on a matrix $X \in \mathbb{R}^{m \times r}$ can be interpreted either as a subset $I \subseteq \llbracket m \rrbracket \times \llbracket r \rrbracket$ or as its indicator matrix $1_I \in \{0, 1\}^{m \times r}$ defined as: $(1_I)_{i, j} = 1$ if $(i, j) \in I$ and 0 otherwise. Both representations will be used interchangeably and the meaning should be clear from the context. For $T \subseteq \llbracket r \rrbracket$, we use the notation $I_T := I \cap (\llbracket m \rrbracket \times T)$ (this is coherent with the notation A_T introduced earlier).

The notation $\text{supp}(A)$ is used for both vectors and matrices: if $A \in \mathbb{R}^m$ is a vector, then $\text{supp}(A) = \{i \mid A_i \neq 0\} \subseteq \llbracket m \rrbracket$; if $A \in \mathbb{R}^{m \times n}$ is a matrix, then $\text{supp}(A) = \{(i, j) \mid A_{i, j} \neq 0\} \subseteq \llbracket m \rrbracket \times \llbracket n \rrbracket$. Given two matrices $A, B \in \mathbb{R}^{m \times n}$, the Hadamard product $A \odot B$ between A and B is defined as $(A \odot B)_{i, j} = A_{i, j} B_{i, j}, \forall (i, j) \in \llbracket m \rrbracket \times \llbracket n \rrbracket$. Since a support constraint I of a matrix X can be thought of as a binary matrix of the same size, we define $X \odot I := X \odot 1_I$ analogously (it is a matrix whose coefficients in I are unchanged while the others are set to zero).

2. Matrix factorization with fixed support is NP-hard. To show that (FSMF) is NP-hard we use the classical technique to prove NP-hardness: reduction. Our choice of reducible problem is matrix completion with noise [5].

DEFINITION 2.1 (Matrix completion with noise [5]). *Let $W \in \{0, 1\}^{m \times n}$ be a binary matrix. Given $A \in \mathbb{R}^{m \times n}$, $s \in \mathbb{N}$, the matrix completion problem (MCP) is:*

$$(MCP) \quad \underset{X \in \mathbb{R}^{m \times s}, Y \in \mathbb{R}^{n \times s}}{\text{Minimize}} \quad \|A - XY^\top\|_W^2 = \|(A - XY^\top) \odot W\|^2.$$

This problem is NP-hard even when $s = 1$ [5] by its reducibility from Maximum-Edge Biclique Problem, which is NP-complete [17]. This is given in the following theorem:

THEOREM 2.2 (NP-hardness of matrix completion with noise [5]). *Given a binary weighting matrix $W \in \{0, 1\}^{m \times n}$ and $A \in [0, 1]^{m \times n}$, the optimization problem*

$$(MCPO) \quad \underset{x \in \mathbb{R}^m, y \in \mathbb{R}^n}{\text{Minimize}} \quad \|A - xy^\top\|_W^2.$$

is called rank-one matrix completion problem (MCPO). Denote p^ the infimum of (MCPO) and let $\epsilon = 2^{-12}(mn)^{-7}$. It is NP-hard to find an approximate solution with objective function accuracy less than ϵ , i.e. with objective value $p \leq p^* + \epsilon$.*

The following lemma gives a reduction from (MCPO) to (FSMF).

LEMMA 2.3. *For any binary matrix $W \in \{0, 1\}^{m \times n}$, there exist an integer r and two sets I and J such that for all $A \in \mathbb{R}^{m \times n}$, (MCPO) and (FSMF) share the same infimum. I and J can be constructed in polynomial time. Moreover, if one of the problems has a known solution that provides objective function accuracy ϵ , we can find a solution with the same accuracy for the other one in polynomial time.*

Proof sketch. Up to a transposition, we can assume without loss of generality that $m \geq n$. Let $r = n + 1 = \min(m, n) + 1$. We define $I \in \{0, 1\}^{m \times (n+1)}$ and

$J \in \{0, 1\}^{n \times (n+1)}$ as follows:

$$I_{i,j} = \begin{cases} 1 - W_{i,j} & \text{if } j \neq n \\ 1 & \text{if } j = n + 1 \end{cases}, J_{i,j} = \begin{cases} 1 & \text{if } j = i \text{ or } j = n + 1 \\ 0 & \text{otherwise} \end{cases}$$

This construction can clearly be made in polynomial time. We show in [Appendix A](#) that the two problems share the same infimum. \square

Using [Lemma 2.3](#), we obtain a result of NP-hardness for [\(FSMF\)](#) as follows.

THEOREM 2.4. *When $A \in [0, 1]^{m \times n}$, it is NP-hard to solve [\(FSMF\)](#) with arbitrary index sets I, J and objective function accuracy less than $\epsilon = 2^{-12}(mn)^{-7}$.*

Proof. Given any instance of [\(MCPO\)](#) (i.e., two matrices $A \in [0, 1]^{m \times n}$ and $W \in \{0, 1\}^{m \times n}$), we can produce an instance of [\(FSMF\)](#) (the same matrix A and $I \in \{0, 1\}^{m \times r}, J \in \{0, 1\}^{n \times r}$) such that both have the same infimum ([Lemma 2.3](#)). Moreover, for any given objective function accuracy, we can use the procedure of [Lemma 2.3](#) to make sure the solutions of both problems share the same accuracy.

Since all procedures are polynomial, this defines a polynomial reduction from [\(MCPO\)](#) to [\(FSMF\)](#). Because [\(MCPO\)](#) is NP-hard to obtain a solution with objective function accuracy less than ϵ ([Theorem 2.2](#)), so is [\(FSMF\)](#). \square

We point out that, while the result is interesting on its own, for some applications, such as those arising in machine learning, the accuracy bound $O((mn)^{-7})$ may not be really appealing. We thus keep as an interesting open research direction to determine if some precision threshold exists that make the general problem easy.

3. Tractable instances of matrix factorization with fixed support. Even though [\(FSMF\)](#) is generally NP-hard, when we consider the full support case $I = \llbracket m \rrbracket \times \llbracket r \rrbracket, J = \llbracket n \rrbracket \times \llbracket r \rrbracket$ (i.e., no coefficients of X, Y are set to zero, they are all optimized), the problem is equivalent to *low rank matrix approximation* (LRMA) [\[3\]](#), which can be solved using the Singular Value Decomposition (SVD) [\[6\]](#)². This section is devoted to enlarge the family of supports for which [\(FSMF\)](#) can be solved by an effective direct algorithm. We start with an important definition:

DEFINITION 3.1 (Support of rank-one contribution). *Given two support constraints $I \in \{0, 1\}^{m \times r}$ and $J \in \{0, 1\}^{n \times r}$ of [\(FSMF\)](#) and $k \in \llbracket r \rrbracket$, we define the k^{th} rank-one contribution support $\mathcal{S}_k(I, J)$ (or in short, \mathcal{S}_k) as: $\mathcal{S}_k(I, J) = I_{\bullet, k} J_{\bullet, k}^\top$. This can be seen either as: a tensor product: $\mathcal{S}_k \in \{0, 1\}^{m \times n}$ is a binary matrix or a Cartesian product: \mathcal{S}_k is a set of matrix indices defined as $\text{supp}(I_{\bullet, k}) \times \text{supp}(J_{\bullet, k})$.*

Given a pair of support constraints I, J , if $\text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$, we have: $\text{supp}(X_{\bullet, k} Y_{\bullet, k}^\top) \subseteq \mathcal{S}_k, \forall k \in \llbracket r \rrbracket$. Since $XY^\top = \sum_{k=1}^r X_{\bullet, k} Y_{\bullet, k}^\top$ the notion of contribution support \mathcal{S}_k captures the constraint on the support of the k^{th} rank-one contribution, $X_{\bullet, k} Y_{\bullet, k}^\top$, of the matrix product XY^\top (illustrated in [Figure 1](#)). We can partition $\llbracket r \rrbracket$ in terms of equivalence classes of rank-one supports:

DEFINITION 3.2 (Equivalence classes of rank-one supports, representative rank-one supports). *Given $I \in \{0, 1\}^{m \times r}, J \in \{0, 1\}^{n \times r}$, define an equivalence relation on $\llbracket r \rrbracket$ as: $i \sim j$ if and only if $\mathcal{S}_i = \mathcal{S}_j$ (or equivalently $(I_{\bullet, i}, J_{\bullet, i}) = (I_{\bullet, j}, J_{\bullet, j})$). This yields a partition of $\llbracket r \rrbracket$ into equivalence classes.*

²Exact SVD is not polynomially tractable, yet it can be practically computed to machine precision in $O(mn^2)$ [\[10\]](#), see also [\[22, Lecture 31, page 236\]](#). It is thus convenient to think of LRMA as polynomially solvable.

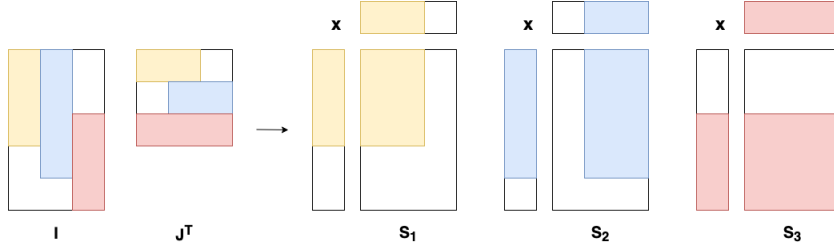


FIG. 1. Illustration the idea of support of rank-one contribution. Colored rectangles indicate the support constraints (I, J) and the support constraints \mathcal{S}_k on each component matrix $X_{\bullet, k} Y_{\bullet, k}^\top$.

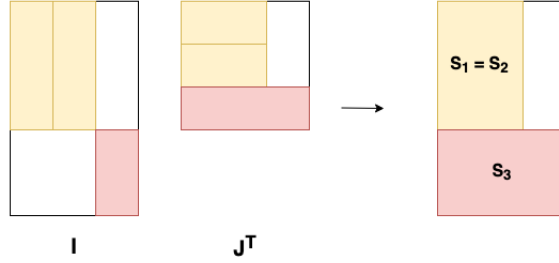


FIG. 2. An instance of support constraints (I, J) satisfying [Theorem 3.3](#). We use colored rectangles to indicate the support constraints (I, J) . The indices belonging to the same equivalence class share the same color.

Denote \mathcal{P} the collection of equivalence classes. For each class $P \in \mathcal{P}$ denote \mathcal{S}_P a representative rank-one support, $R_P \subseteq \llbracket m \rrbracket$ and $C_P \subseteq \llbracket n \rrbracket$ the supports of rows and columns in \mathcal{S}_P , respectively. For every $k \in P$ we have $\mathcal{S}_k = \mathcal{S}_P$ and $\text{supp}(I_{\bullet, k}) = R_P$, $\text{supp}(J_{\bullet, k}) = C_P$.

For every $\mathcal{P}' \subseteq \mathcal{P}$ denote $\mathcal{S}_{\mathcal{P}'} = \cup_{P \in \mathcal{P}'} \mathcal{S}_P \subseteq \llbracket m \rrbracket \times \llbracket n \rrbracket$ and $\bar{\mathcal{S}}_{\mathcal{P}'} = (\llbracket m \rrbracket \times \llbracket n \rrbracket) \setminus \mathcal{S}_{\mathcal{P}'}$.

For instance, in the example in [Figure 1](#) we have three distinct equivalent classes. A first simple sufficient condition ensuring the tractability of an instance of [\(FSMF\)](#) is as follows.

THEOREM 3.3. Consider $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$, and \mathcal{P} the collection of equivalence classes of [Definition 3.2](#). If the representative rank-one supports are pairwise disjoint, i.e., $\mathcal{S}_P \cap \mathcal{S}_{P'} = \emptyset$ for each distinct $P, P' \in \mathcal{P}$, then matrix factorization with fixed support is tractable for any $A \in \mathbb{R}^{m \times n}$.

Proof. In this proof, for each equivalent class $P \in \mathcal{P}$ ([Definition 3.2](#)) we use the notations $X_P \in \mathbb{R}^{m \times r}$, $Y_P \in \mathbb{R}^{n \times r}$ (introduced in [Subsection 1.1](#)). We also use the notations R_P, C_P ([Definition 3.2](#)). For each equivalent class P , we have:

$$(3.1) \quad (X_P Y_P^\top)_{R_P, C_P} = X_{R_P, P} Y_{C_P, P}^\top$$

and the product XY^\top can be decomposed as: $XY^\top = \sum_{P \in \mathcal{P}} X_P Y_P^\top$. Due to the hypothesis of this theorem, with $P, P' \in \mathcal{P}$, $P' \neq P$, we further have:

$$(3.2) \quad X_{P'} Y_{P'}^\top \odot \mathcal{S}_P = \mathbf{0}$$

Algorithm 3.1 Fixed support matrix factorization (under [Theorem 3.3](#) assumptions)

- 1: **procedure** SVD_FSMF($A \in \mathbb{R}^{m \times n}$, $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$)
 - 2: Initialize $X = \mathbf{0}$, $Y = \mathbf{0}$.
 - 3: Partition $\llbracket r \rrbracket$ into \mathcal{P} ([Definition 3.2](#)).
 - 4: **for** $P \in \mathcal{P}$ **do**
 - 5: Compute the truncated SVD of A_{R_P, C_P} to find a pair (X^*, Y^*) , $X^* \in \mathbb{R}^{|R_P| \times |P|}$, $Y^* \in \mathbb{R}^{|C_P| \times |P|}$ that minimizes $\|A_{R_P, C_P} - X^*(Y^*)^\top\|^2$.
 - 6: Assign $X_{R_P, P} = X^*$, $Y_{C_P, P} = Y^*$.
 - 7: **end for**
 - 8: **return** (X, Y)
 - 9: **end procedure**
-

The objective function $L(X, Y)$ is:

$$\begin{aligned}
 (3.3) \quad \|A - XY^\top\|^2 &= \left(\sum_{P \in \mathcal{P}} \|(A - XY^\top) \odot \mathcal{S}_P\|^2 \right) + \|(A - XY^\top) \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\
 &= \left(\sum_{P \in \mathcal{P}} \|(A - \sum_{P' \in \mathcal{P}} X_{P'} Y_{P'}^\top) \odot \mathcal{S}_P\|^2 \right) + \|(A - \sum_{P' \in \mathcal{P}} X_{P'} Y_{P'}^\top) \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\
 &\stackrel{(3.2)}{=} \left(\sum_{P \in \mathcal{P}} \|(A - X_P Y_P^\top) \odot \mathcal{S}_P\|^2 \right) + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\
 &= \left(\sum_{P \in \mathcal{P}} \|A_{R_P, C_P} - (X_P Y_P^\top)_{R_P, C_P}\|^2 \right) + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\
 &\stackrel{(3.1)}{=} \left(\sum_{P \in \mathcal{P}} \|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2 \right) + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2
 \end{aligned}$$

Therefore, if we ignore the constant term $\|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2$, the function $L(X, Y)$ is decomposed into a sum of functions $\|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2$, which are LRMA instances. Since all the optimized parameters are $\{(X_{R_P, P}, Y_{C_P, P})\}_{P \in \mathcal{P}}$, an optimal solution of L is $\{(X_{R_P, P}^*, Y_{C_P, P}^*)\}_{P \in \mathcal{P}}$, where $(X_{R_P, P}^*, Y_{C_P, P}^*)$ is a minimizer of $\|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2$. Since $(X_{R_P, P}^*, Y_{C_P, P}^*)$ can be calculated by SVD, the problem can be solved efficiently. \square

For these easy instances, we can therefore recover the factors in polynomial time with the procedure described in [Algorithm 3.1](#). Given a target matrix $A \in \mathbb{R}^{m \times n}$ and support constraints $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$ satisfying the condition in [Theorem 3.3](#), [Algorithm 3.1](#) returns two factors (X, Y) solution of (FSMF).

[Theorem 3.3](#) requires all the rank-one contribution supports of different equivalence classes to be disjoint (as illustrated on [Figure 2](#)). Although this assumption appears restrictive, it is verified for certain interesting support constraints in practice. In [\[11\]](#), we show such an example. We also propose in [\[11\]](#) a hierarchical extension of our method, designed to handle multi-layer matrix factorization (the case in which the matrix is approximated as the product of more than two factors) and demonstrate the superior performance of [Algorithm 3.3](#) in comparison to first-order optimization approaches commonly used for the solution of such problems, in terms both of computational time and accuracy.

In the next result, we explore the tractability of (FSMF) while allowing partial intersection between two representative rank-one contribution supports.

DEFINITION 3.4 (Complete equivalence classes of rank-one supports - CEC). *$P \in \mathcal{P}$ is a complete equivalence class (or CEC) if $|P| \geq \min\{|C_P|, |R_P|\}$ with C_P, R_P as in [Definition 3.2](#). Denote $\mathcal{P}^* \subseteq \mathcal{P}$ the family of all complete equivalence classes, $T = \cup_{P \in \mathcal{P}^*} P \subseteq \llbracket r \rrbracket$, $\bar{T} = \llbracket r \rrbracket \setminus T$, and the shorthand $\mathcal{S}_T = \mathcal{S}_{\mathcal{P}^*}$.*

The interest of complete equivalence classes is that their expressivity is powerful enough to represent any matrix whose support is included in \mathcal{S}_T , as illustrated by the following lemma associated to [Algorithm 3.2](#).

LEMMA 3.5. *Given $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$, consider T, \mathcal{S}_T as in [Definition 3.4](#). For any matrix $A \in \mathbb{R}^{m \times n}$ such that $\text{supp}(A) \subseteq \mathcal{S}_T$, there exist $X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}$ such that $A = XY^\top$ and $\text{supp}(X) \subseteq I_T, \text{supp}(Y) \subseteq J_T$. Such a pair can be computed using [Algorithm 3.2](#) with input A, I, J .*

Algorithm 3.2 Find (X, Y) satisfying [Lemma 3.5](#)'s assumptions

```

1: procedure FILL_CEC( $A \in \mathbb{R}^{m \times n}, I \in \{0, 1\}^{m \times r}, J \in \{0, 1\}^{n \times r}$ )
2:   Initialize  $X = \mathbf{0}, Y = \mathbf{0}$ .
3:   Compute  $\mathcal{P}^*$  from  $(I, J)$  (Definition 3.4).
4:   for  $P \in \mathcal{P}^* := \{P_1, \dots, P_\ell\}$  do
5:     Let  $A' = A - XY^\top$ 
6:     if  $|P| \geq |R_P|$  then
7:       Choose an arbitrary matrix  $X' \in \mathbb{R}^{|R_P| \times |P|}$  with full row rank.
8:       Assign  $X_{R_P, P} = X', Y_{C_P, P} = (X'^\top (X' X'^\top)^{-1} A'_{R_P, C_P})^\top$ .
9:     else  $\triangleright$  Here necessarily  $|P| \geq |C_P|$ , since  $|P| \geq \min\{|C_P|, |R_P|\}$ 
10:      Choose an arbitrary matrix  $Y' \in \mathbb{R}^{|C_P| \times |P|}$  with full row rank.
11:      Assign  $X_{R_P, P} = A'_{R_P, C_P} (Y' Y'^\top)^{-1} Y', Y_{C_P, P} = Y'$ .
12:     end if
13:   end for
14:   return  $(X, Y)$ 
15: end procedure

```

The proof of [Lemma 3.5](#) is deferred to the supplementary material ([Appendix B.1](#)). The next definition introduces the key properties that the indices $k \in \llbracket r \rrbracket$ which are not in any CEC need to satisfy in order to make (FSMF) overall tractable.

DEFINITION 3.6 (Rectangular support outside CECs of rank-one supports). *Given $I \in \{0, 1\}^{m \times r}, J \in \{0, 1\}^{n \times r}$, consider T and \mathcal{S}_T as in [Definition 3.4](#) and $\bar{T} = \llbracket r \rrbracket \setminus T$. For $k \in \bar{T}$ define the support outside CECs of the k^{th} rank-one support. as: $\mathcal{S}'_k = \mathcal{S}_k \setminus \mathcal{S}_T$. If $\mathcal{S}'_k = R_k \times C_k$ for some $R_k \subseteq \llbracket m \rrbracket, C_k \subseteq \llbracket n \rrbracket$, (or equivalently \mathcal{S}'_k is of rank at most one), we say the support outside CECs of the k^{th} rank-one support \mathcal{S}'_k is rectangular.*

To state our tractability result, we further categorize the indices in I and J as follows:

DEFINITION 3.7 (Taxonomy of indices of I and J). *With the notations of [Definition 3.6](#), assume that \mathcal{S}'_k is rectangular for all $k \in \bar{T}$. We decompose the indices of I (resp J) into three sets as follows:*

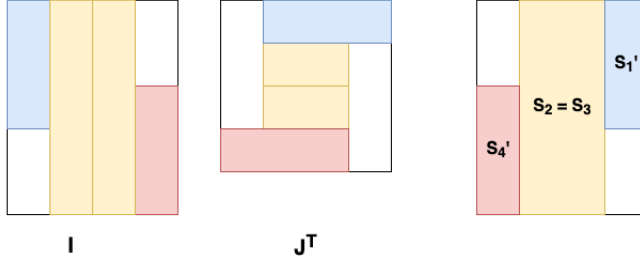


FIG. 3. An instance of support constraints (I, J) satisfying the assumptions of [Theorem 3.8](#). We have $T = \{2, 3\}$. The supports outside CEC S'_1 and S'_4 are disjoint.

	Classification for I	Classification for J
1	$I_T = \{(i, k) \mid k \in T, i \in \llbracket m \rrbracket\} \cap I$	$J_T = \{(j, k) \mid k \in T, j \in \llbracket n \rrbracket\} \cap J$
2	$I_T^1 = \{(i, k) \mid k \notin T, i \in R_k\} \cap I$	$J_T^1 = \{(j, k) \mid k \notin T, j \in C_k\} \cap J$
3	$I_T^2 = \{(i, k) \mid k \notin T, i \notin R_k\} \cap I$	$J_T^2 = \{(j, k) \mid k \notin T, j \notin C_k\} \cap J$

The following theorem generalizes [Theorem 3.3](#).

THEOREM 3.8. Consider $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$. Assume that for all $k \in \bar{T}$, S'_k is rectangular and that for all $k, l \in \bar{T}$ we have $S'_k = S'_l$ or $S'_k \cap S'_l = \emptyset$. Then, $(I_{\bar{T}}^1, J_{\bar{T}}^1)$ satisfy the assumptions of [Theorem 3.3](#). Moreover, for any matrix $A \in \mathbb{R}^{m \times n}$, two instances of (FSMF) with data (A, I, J) and $(A \odot \bar{S}_T, I_{\bar{T}}^1, J_{\bar{T}}^1)$ respectively, share the same infimum. Given an optimal solution of one instance, we can construct the optimal solution of the other in polynomial time.

[Theorem 3.8](#) is proved in [Appendix B.2](#). It implies that solving the problem with support constraints (I, J) can be achieved by reducing to another problem, with support constraints satisfying the assumptions of [Theorem 3.3](#). The latter problem can thus be efficiently solved by [Algorithm 3.1](#). In particular, [Theorem 3.3](#) is a special case of [Theorem 3.8](#) when all the equivalent classes (including CECs) have disjoint representative rank-one supports.

[Figure 3](#) shows an instance of (I, J) satisfying the assumptions of [Theorem 3.8](#). An algorithm for instances satisfying the assumptions of [Theorem 3.8](#) is given in [Algorithm 3.3](#) (more details can be found in [Corollary B.3](#) and [Remark B.4](#) in [Appendix B](#)). In [Algorithm 3.3](#), [Algorithm 3.1](#) is used at Line 3, and [Algorithm 3.2](#) at Line 4.

Algorithm 3.3 Fixed support matrix factorization (under [Theorem 3.8](#)'s assumptions)

- 1: **procedure** SVD_FSMF2($A \in \mathbb{R}^{m \times n}, I \in \{0, 1\}^{m \times r}, J \in \{0, 1\}^{n \times r}$)
 - 2: Partition the indices of I, J into I_T, I_T^1, I_T^2 (and J_T, J_T^1, J_T^2) ([Definition 3.6](#)).
 - 3: $(X_{\bar{T}}^1, Y_{\bar{T}}^1) = \text{SVD_FSMF}(A \odot \bar{S}_T, I_{\bar{T}}^1, J_{\bar{T}}^1)$ (T, \bar{S}_T as in [Definition 3.4](#)).
 - 4: $(X_T, Y_T) = \text{FILL_CEC}(A \odot \bar{S}_T, I, J)$.
 - 5: **return** $(X_T + X_{\bar{T}}^1, Y_T + Y_{\bar{T}}^1)$
 - 6: **end procedure**
-

4. Landscape of matrix factorization with fixed support. In this section, we first recall the definition of *spurious local valleys* and *spurious local minima*, which

are undesirable objects in the landscape of a function, as they may prevent local optimization methods to converge to globally optimal solutions. Previous works [24, 25, 9] showed that the landscape of the optimization problem associated to low rank approximation is free of such *spurious objects*, which potentially gives the intuition for its tractability.

We prove that similar results hold for the much richer family of tractable support constraints for (FSMF) that we introduced in Theorem 3.3. The landscape with the assumptions of Theorem 3.8 is also analyzed. These results might suggest a natural conjecture: an instance of (FSMF) is tractable if and only if the landscape is benign. However, this is not true. We show an example that contradicts this conjecture: we show an instance of (FSMF) that can be solved efficiently, despite the fact that its corresponding landscape contains spurious objects. We will see in the next section that the opposite direction is not so evident either: we propose a numerical illustration of the fact that even when the landscape is benign, the solution of (FSMF) may not be so straightforward with standard iterative methods.

4.1. Spurious local minima and spurious local valleys. We start by recalling the classical definitions of global and local minima of a real-valued function.

DEFINITION 4.1 (Spurious local minimum [25, 16]). *Consider $L : \mathbb{R}^d \rightarrow \mathbb{R}$. A vector $x^* \in \mathbb{R}^d$ is a:*

- **global minimum** (of L) if $L(x^*) \leq L(x), \forall x$.
- **local minimum** if there is a neighborhood \mathcal{N} of x^* such that $L(x^*) \leq L(x), \forall x \in \mathcal{N}$.
- **strict local minimum** if there is a neighborhood \mathcal{N} of x^* such that $L(x^*) < L(x), \forall x \in \mathcal{N}, x \neq x^*$.
- **(strict) spurious local minimum** if x^* is a (strict) local minimum but it is not a global minimum.

The presence of spurious local minima is undesirable because local optimization methods can get stuck in one of them and never reach the global optimum.

Remark 4.2. With the loss functions $L(X, Y)$ considered in this paper, strict local minima do not exist since for every invertible diagonal matrix D , possibly arbitrarily close to the identity, we have $L(XD, YD^{-1}) = L(X, Y)$.

However, this is not the only undesirable landscape in an optimization problem: spurious local valleys, as defined next, are also challenging.

DEFINITION 4.3 (Sublevel Set [1]). *Consider $L : \mathbb{R}^d \rightarrow \mathbb{R}$. For every $\alpha \in \mathbb{R}$, the α -level set of L is the set $E_\alpha = \{x \in \mathbb{R}^d \mid L(x) \leq \alpha\}$.*

DEFINITION 4.4 (Path-Connected Set and Path-Connected Component). *A subset $S \subseteq \mathbb{R}^d$ is path-connected if for every $x, y \in S$, there is a continuous function $r : [0, 1] \rightarrow S$ such that $r(0) = x, r(1) = y$. A path-connected component of $E \subseteq \mathbb{R}^d$ is a maximal path-connected subset: $S \subseteq E$ is path-connected, and if $S' \subseteq E$ is path-connected with $S \subseteq S'$ then $S = S'$.*

DEFINITION 4.5 (Spurious Local Valley [24, 15]). *Consider $L : \mathbb{R}^d \rightarrow \mathbb{R}$ and a set $S \subset \mathbb{R}^d$.*

- S is a **local valley** of L if it is a non-empty path-connected component of some sublevel set.
- S is a **spurious local valley** of L if it is a local valley of L and does not contain a global minimum.

The notion of spurious local valley is inspired by the definition of a *strict* spurious local minimum. If x^* is a strict spurious local minimum, then $\{x^*\}$ is a spurious local

valley. However, the notion of spurious local valley has a wider meaning than just a neighborhood of a strict spurious local minimum. Figure 4 illustrates some other scenarios: as shown on Figure 4a, the segment (approximately) $[10, +\infty)$ creates a

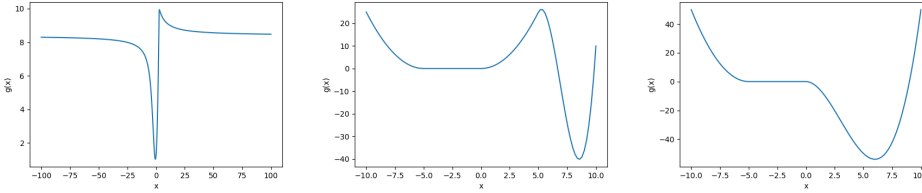


FIG. 4. Examples of functions with spurious objects.

spurious local valley, and this function has only one local (and global) minimizer, at zero; in Figure 4b, there are spurious local minima that are not strict, but form a spurious local valley anyway. It is worth noticing that the concept of a spurious local valley does *not* cover that of a spurious local minimum. Functions can have spurious (non-strict) local minima even if they do not possess any spurious local valley (Figure 4c). Therefore, in this paper, we treat the existence of spurious local valleys and spurious local minima independently. The common point is that if the landscape possesses either of them, local optimization methods need to have proper initialization to have guarantees of convergence to a global minimum.

4.2. Previous results on the landscape. Previous works [9, 25] studied the non-existence of spurious local minima of (FSMF) in the classical case of “low rank matrix approximation” (or *full support matrix factorization*)³. To prove that a critical point is never a spurious local minimum, previous work used the notion of *strict saddle point* (i.e a point where the Hessian is not positive semi-definite, or equivalently has at least one *strictly* – negative eigenvalue), see Definition 4.10 below. To prove the non-existence of spurious local valleys, the following lemma was employed in previous works [24, 15]:

LEMMA 4.6 (Sufficient condition for the non-existence of any spurious local valley [24, Lemma 2]). *Consider a continuous function $L : \mathbb{R}^d \rightarrow \mathbb{R}$. Assume that, for any initial parameter $\tilde{x} \in \mathbb{R}^d$, there exists a continuous path $f : t \in [0, 1] \rightarrow \mathbb{R}^d$ such that:*

- a) $f(0) = \tilde{x}$.
- b) $f(1) \in \arg \min_{x \in \mathbb{R}^d} L(x)$.
- c) *The function $L \circ f : t \in [0, 1] \rightarrow \mathbb{R}$ is non-increasing.*

Then there is no spurious local valley in the landscape of function L .

The result is intuitive and a formal proof can be found in [24]. The theorem claims that given any initial point, if one can find a continuous path connecting the initial point to a global minimizer and the loss function is non-increasing on the path, then there does not exist any spurious local valley. We remark that although (FSMF) is a constrained optimization problem, Lemma 4.6 is still applicable because one can think of the objective function as defined on a subspace: $L : \mathbb{R}^{|I|+|J|} \rightarrow \mathbb{R}$. In this work, to apply Lemma 4.6, the constructed function f has to be a *feasible path*, defined as:

³Since previous works also considered the case $r \geq m, n$, low rank approximation might be misleading sometimes. That is why we occasionally use the name full support matrix factorization to emphasize this fact., where no support constraints are imposed ($I = \llbracket m \rrbracket \times \llbracket r \rrbracket, J = \llbracket n \rrbracket \times \llbracket r \rrbracket$)

DEFINITION 4.7 (Feasible path). *A feasible path w.r.t the support constraints (I, J) (or simply a feasible path) is a continuous function $f(t) = (X_f(t), Y_f(t)) : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$ satisfying $\text{supp}(X_f(t)) \subseteq I, \text{supp}(Y_f(t)) \subseteq J, \forall t \in [0, 1]$.*

Conversely, we generalize and formalize an idea from [24] into the following lemma, which gives a sufficient condition for the existence of a spurious local valley:

LEMMA 4.8 (Sufficient condition for the existence of a spurious local valley). *Consider a continuous function $L : \mathbb{R}^d \rightarrow \mathbb{R}$ whose global minimum is attained. Assume we know three subsets $S_1, S_2, S_3 \subset \mathbb{R}^d$ such that:*

- 1) *The global minima of L are in S_1 .*
- 2) *Every continuous path from S_3 to S_1 passes through S_2 .*
- 3) $\inf_{x \in S_2} L(x) > \inf_{x \in S_3} L(x) > \inf_{x \in S_1} L(x)$.

Then L has a spurious local valley.

Proof. Denote $\Sigma = \{x \mid L(x) = \inf_{x \in \mathbb{R}^d} L(\theta)\}$ the set of global minimizers of L . Σ is not empty due to the assumption that the global minimum is attained, and $\Sigma \subseteq S_1$ by the first assumption.

Since $\inf_{x \in S_2} L(x) > \inf_{x \in S_3} L(x)$, there exists $\tau \in S_3, L(\tau) < \inf_{x \in S_2} L(x)$. Consider Φ the path-connected component of the sublevel set $\{x \mid L(x) \leq L(\tau)\}$ that contains τ . Since Φ is a non-empty path-connected component of a level set, it is a local valley. It is thus sufficient to prove that $\Phi \cap \Sigma = \emptyset$ to obtain that it matches the very definition of a spurious local valley.

Indeed, by contradiction, let's assume that there exists $\tau' \in \Phi \cap \Sigma$. Since $\tau, \tau' \in \Phi$ and Φ is path-connected, by definition of path-connectedness there exists a continuous function $f : [0, 1] \rightarrow \Phi$ such that $f(0) = \tau \in S_3, f(1) = \tau' \in \Sigma \subseteq S_1$. Due to the assumption that every continuous path from S_3 to S_1 has to pass through a point in S_2 , there must exist $t \in (0, 1)$ such that $f(t) \in S_2 \cap \Phi$. Therefore, $L(f(t)) \leq L(\tau)$ (since $f(t) \in \Phi$) and $L(f(t)) > L(\tau)$ (since $f(t) \in S_2$), which is a contradiction. \square

To finish this section, we formally recall previous results which are related to (FSMF) and will be used in our subsequent proofs. The questions of the existence of spurious local valleys and spurious local minima were addressed in previous works for full support matrix factorization and deep linear neural networks [24, 15, 25, 9]. We present only results related to our problem of interest.

THEOREM 4.9 (No spurious local valleys in linear networks [24, Theorem 11]). *Consider linear neural networks of any depth $K \geq 1$ and of any layer widths $p_k \geq 1$ and any input - output dimension $n, m \geq 1$ with the following form: $\Phi(b, \theta) = W_K \dots W_1 b$ where $\theta = (W_i)_{i=1}^K$, and $b \in \mathbb{R}^n$ is a training input sample. With the squared loss function, there is no spurious local valley. More specifically, the function $L(\theta) = \|A - \Phi(B, \theta)\|^2$ satisfies the condition of Lemma 4.6 for any matrices $A \in \mathbb{R}^{m \times N}$ and $B \in \mathbb{R}^{n \times N}$ (A and B are the whole sets of training output and input respectively).*

DEFINITION 4.10 (Strict saddle point property [25, Definition 3]). *Consider a twice differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$. If each critical point of f is either a global minimum or a strict saddle point then f is said to have the strict saddle point property. When this property holds, f has no spurious local minimum.*

Even if f has the strict saddle point property, it may have no global minimum, consider e.g. the function $f(x) = -\|x\|_2^2$.

THEOREM 4.11 (No spurious local minima in shallow linear networks [25, Theorem 3]). *Let $B \in \mathbb{R}^{d_0 \times N}, A \in \mathbb{R}^{d_2 \times N}$ be input and output training examples. Consider the*

problem:

$$\text{Minimize}_{X \in \mathbb{R}^{d_0 \times d_1}, Y \in \mathbb{R}^{d_1 \times d_2}} L(X, Y) = \|A - XYB\|^2$$

If B is full row rank, f has the strict saddle point property (see [Definition 4.10](#)) hence f has no spurious local minimum.

Both theorems are valid for a particular case of matrix factorization with fixed support: full support matrix factorization. Indeed, given a factorized matrix $A \in \mathbb{R}^{m \times n}$, in [Theorem 4.9](#), if $K = 2, B = \mathbf{I}_n$ ($n = N$), then the considered function is $L = \|A - W_2 W_1\|^2$. This is (FSMF) without support constraints I and J (and without a transpose on W_1 , which does not change the nature of the problem). [Theorem 4.9](#) guarantees that L satisfies the conditions of [Lemma 4.6](#), thus has no spurious local valley.

Similarly, in [Theorem 4.11](#), if $B = \mathbf{I}_{d_0}$ ($d_0 = N$, therefore B is full row rank), we return to the same situation of [Theorem 4.9](#). In general, [Theorem 4.11](#) claims that the landscape of the full support matrix factorization problem has the strict saddle point property and thus, does not have spurious local minima.

However, once we turn to (FSMF) with *arbitrary* I and J , such benign landscape is not guaranteed anymore, as we will show in [Example 4.22](#). Our work in the next subsections studies conditions on the support constraints I and J ensuring the absence / allowing the presence of spurious objects, and can be considered as a generalization of previous results with full supports. [[25](#), [24](#), [9](#)].

4.3. Landscape of matrix factorization with fixed support constraints.

We start with the first result on the landscape in the simple setting of [Theorem 3.3](#).

THEOREM 4.12. *Under the assumption of [Theorem 3.3](#), the function $L(X, Y)$ in (FSMF) does not admit any spurious local valley for any matrix A . In addition, L has the strict saddle point property.*

Proof. Recall that under the assumption of [Theorem 3.3](#), all the variables to be optimized are decoupled into “blocks” $\{(X_{R_P, P}, Y_{C_P, P})\}_{P \in \mathcal{P}}$ (P, \mathcal{P} are defined in [Definition 3.2](#)). We denote $\mathcal{P} = \{P_1, P_2, \dots, P_\ell\}$, $P_i \subseteq \llbracket r \rrbracket$, $1 \leq i \leq \ell$. From [Equation \(3.3\)](#), we have:

$$(4.1) \quad \|A - XY^\top\|^2 = \left(\sum_{P \in \mathcal{P}} \|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2 \right) + \|A \odot \bar{\mathcal{S}}\|^2$$

Therefore, the function $L(X, Y)$ is a sum of functions $L_P(X_{R_P, P}, Y_{C_P, P}) := \|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2$, which do *not* share parameters and are instances of the full support matrix factorization problem restricted to the corresponding blocks in A . The global minimizers of L are $\{(X_{R_P, P}^*, Y_{C_P, P}^*)\}_{P \in \mathcal{P}}$, where for each $P \in \mathcal{P}$ the pair $(X_{R_P, P}^*, Y_{C_P, P}^*)$ is any global minimizer of $\|A_{R_P, C_P} - X_{R_P, P} Y_{C_P, P}^\top\|^2$.

1) **Non-existence of any spurious local valley:** By [Theorem 4.9](#), from any initial point $(X_{R_P, P}^0, Y_{C_P, P}^0)$, there exists a continuous function $f_P(t) = (\tilde{X}_P(t), \tilde{Y}_P(t)) : [0, 1] \mapsto \mathbb{R}^{|R_P| \times |P|} \times \mathbb{R}^{|C_P| \times |P|}$ satisfying the conditions in [Lemma 4.6](#), which are:

- i) $f_P(0) = (X_{R_P, P}^0, Y_{C_P, P}^0)$.
- ii) $f_P(1) = (X_{R_P, P}^*, Y_{C_P, P}^*)$.
- iii) $L_P \circ f_P : [0, 1] \rightarrow \mathbb{R}$ is non-increasing.

Consider a feasible path ([Definition 4.7](#)) $f(t) = (\tilde{X}(t), \tilde{Y}(t)) : [0, 1] \mapsto \mathbb{R}^{m \times r} \times \mathbb{R}^{r \times n}$ defined in such a way that $\tilde{X}(t)_{R_P, P} = \tilde{X}_P(t)$ for each $P \in \mathcal{P}$ and similarly for

$\tilde{Y}(t)$. Since $L \circ f = \sum_{P \in \mathcal{P}} L_P \circ f_P + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2$, f satisfies the assumptions of [Lemma 4.6](#), which shows the non-existence of any spurious local valley.

- 2) **Non-existence of any spurious local minimum:** Due to the decomposition in [Equation \(4.1\)](#), the gradient and Hessian of $L(X, Y)$ have the following form:

$$\frac{\partial L}{\partial X_{R_P, P}} = \frac{\partial L_P}{\partial X_{R_P, P}}, \quad \frac{\partial L}{\partial Y_{C_P, P}} = \frac{\partial L_P}{\partial Y_{C_P, P}}, \quad \forall P \in \mathcal{P}$$

$$H(L)_{|(X, Y)} \begin{pmatrix} H(L_{P_1})_{|(X_{R_{P_1}, P_1}, Y_{C_{P_1}, P_1})} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & H(L_{P_\ell})_{|(X_{R_{P_\ell}, P_\ell}, Y_{C_{P_\ell}, P_\ell})} \end{pmatrix}$$

Consider a critical point (X, Y) of $L(X, Y)$ that is not a global minimizer. Since (X, Y) is a critical point of $L(X, Y)$, $(X_{R_P, P}, Y_{C_P, P})$ is a critical point of the function L_P for all $P \in \mathcal{P}$. Since (X, Y) is not a global minimizer of $L(X, Y)$, there exists $P \in \mathcal{P}$ such that $(X_{R_P, P}, Y_{C_P, P})$ is not a global minimizer of L_P . By [Theorem 4.11](#), $H(L_P)_{|(X_{R_P, P}, Y_{C_P, P})}$ is not positive semi-definite. Hence, $H(L)_{|(X, Y)}$ is not positive semi-definite either (since $H(L)_{|(X, Y)}$ has block diagonal form). This implies that (X, Y) is a strict saddle point as well (hence, not a spurious local minimum). \square

For spurious local valleys, we have the same results for the setting in [Theorem 3.8](#). The proof is, however, less straightforward.

THEOREM 4.13. *If I, J satisfy the assumptions of [Theorem 3.8](#), then for each matrix A the landscape of $L(X, Y)$ in (FSMF) has no spurious local valley.*

The following is a concept which will be convenient for the proof of [Theorem 4.13](#).

DEFINITION 4.14 (CEC-full-rank). *A feasible point (X, Y) is said CEC-full-rank if $\forall P \in \mathcal{P}^*$, either $X_{R_P, P}$ or $Y_{C_P, P}$ is full row rank.*

We need three following lemmas to prove [Theorem 4.13](#):

LEMMA 4.15. *Given $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$, consider T and \mathcal{S}_T as in [Definition 3.2](#) and a feasible point (X, Y) . There exists a feasible path $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

- 1) f connects (X, Y) with a CEC-full-rank point: $f(0) = (X, Y)$, and $f(1)$ is CEC-full-rank.
- 2) $X_f(t)(Y_f(t))^\top = XY^\top, \forall t \in [0, 1]$.

LEMMA 4.16. *Under the assumption of [Theorem 3.8](#), for any CEC-full-rank feasible point (X, Y) , there exists feasible path $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

- 1) $f(0) = (X, Y)$.
- 2) $L \circ f$ is non-increasing.
- 3) $(A - X_f(1)(Y_f(1))^\top) \odot \mathcal{S}_T = \mathbf{0}$.

LEMMA 4.17. *Under the assumption of [Theorem 3.8](#), for any CEC-full-rank feasible point (X, Y) verifying: $(A - XY^\top) \odot \mathcal{S}_T = \mathbf{0}$, there exists a feasible path $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

- 1) $f(0) = (X, Y)$.
- 2) $L \circ f$ is non-increasing.
- 3) $f(1)$ is an optimal solution of L .

The proofs of [Lemma 4.15](#), [Lemma 4.16](#) and [Lemma 4.17](#) can be found in [Appendix D.1](#), [Appendix D.2](#) and [Appendix D.3](#).

Proof of Theorem 4.13. Given any initial point (X^0, Y^0) , Lemma 4.15 shows the existence of a continuous path along which the product of $XY^\top = X^0(Y^0)^\top$ does not change (thus, $L(X, Y)$ is constant) and ending at a CEC-full-rank point. Therefore it is sufficient to prove the theorem under the additional assumption that (X^0, Y^0) is CEC-full-rank. With this additional assumption, one can employ Lemma 4.16 to build a continuous path $f_1(t) = (X_1(t), Y_1(t))$, such that $t \mapsto L(X_1(t), Y_1(t))$ is non-increasing, that connects (X^0, Y^0) to a point (X^1, Y^1) satisfying:

$$(A - X^1(Y^1)^\top) \odot \mathcal{S}_T = \mathbf{0}.$$

Again, one can assume that (X^1, Y^1) is CEC-full-rank (one can invoke Lemma 4.15 one more time). Therefore, (X^1, Y^1) satisfies the conditions of Lemma 4.17. Hence, there exists a continuous path $f_2(t) = (X_2(t), Y_2(t))$ that makes $L(X_2(t), Y_2(t))$ non-increasing and that connects (X^1, Y^1) to (X^*, Y^*) , a global minimizer.

Finally, since the concatenation of f_1 and f_2 satisfies the assumptions of Lemma 4.6, we can conclude that there is no spurious local valley in the landscape of $\|A - XY^\top\|^2$. \square

The next natural question is whether spurious local minima exist in the setting of Theorem 3.8. While in the setting of Theorem 3.3, all critical points which are not global minima are saddle points, the setting of Theorem 3.8 allows second order critical points (point whose gradient is zero and Hessian is positive semi-definite), which are not global minima.

Example 4.18. Consider the following pair of support constraints I, J and factorized matrix $I = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, $J = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$, $A = \begin{bmatrix} 10 & 0 \\ 0 & 1 \end{bmatrix}$. With the notations of Definition 3.4 we have $T = \{1\}$ and one can check that this choice of I and J satisfies the assumptions of Theorem 3.8. The infimum of $L(X, Y) = \|A - XY^\top\|^2$ is zero, and attained, for example at $X^* = I_2, Y^* = A$. Consider the following feasible point (X_0, Y_0) : $X_0 = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$, $Y_0 = \begin{bmatrix} 0 & 10 \\ 0 & 0 \end{bmatrix}$. Since $X_0 Y_0^\top = \begin{bmatrix} 10 & 0 \\ 0 & 0 \end{bmatrix} \neq A$, (X_0, Y_0) is not a global optimal solution. Calculating the gradient of L verifies that (X_0, Y_0) is a critical point:

$$\nabla L(X_0, Y_0) = ((A - X_0 Y_0^\top) Y_0, (A^\top - Y_0 X_0^\top) X_0) = (\mathbf{0}, \mathbf{0})$$

Nevertheless, the Hessian of the function L at (X_0, Y_0) is positive semi-definite. Direct calculation can be found in Appendix D.5.

This example shows that if we want to prove the non-existence of spurious local minima in the new setting, one cannot rely on the Hessian. This is challenging since the second order derivatives computation is already tedious. Nevertheless, with Definition 4.14, we can still say something about spurious local minima in the new setting.

THEOREM 4.19. *Under the assumptions of Theorem 3.8, if a feasible point (X, Y) is CEC-full-rank, then (X, Y) is not a spurious local minimum of (FSMF). Otherwise there is a feasible path, along which $L(\cdot, \cdot)$ is constant, that joins (X, Y) to some (\tilde{X}, \tilde{Y}) which is not a spurious local minimum.*

When (X, Y) is not CEC-full-rank, the theorem guarantees that it is not a strict local minimum, since there is path starting from (X, Y) with constant loss. This should however not be a surprise in light of Remark 4.2: indeed, the considered loss function admits no strict local minimum at all. Yet, the path with “flat” loss constructed in the theorem is fundamentally different from the ones naturally due to scale invariances of the problem and captured by Remark 4.2. Further work would be needed to investigate whether this can be used to get a stronger result.

Proof sketch. To prove this theorem, we proceed through two main steps:

- 1) First, we show that any local minimum satisfies:

$$(4.2) \quad (A - XY^\top) \odot \mathcal{S}_T = \mathbf{0}$$

- 2) Second, we show that if a point (X, Y) is CEC-full-rank and satisfies Equation (4.2), it cannot be a spurious local minimum.

Combining the above to steps, we obtain as claimed that if a feasible pair (X, Y) is CEC-full-rank, then it is not a spurious local minimum. Finally, if a feasible pair (X, Y) is not CEC-full-rank, Lemma 4.15 yields a feasible path along which L is constant that joins (X, Y) to some feasible (\tilde{X}, \tilde{Y}) which is CEC-full-rank, hence (as we have just shown) not a spurious local minimum.

A complete proof is presented in Appendix D.4. \square

Although Theorem 4.19 does not exclude completely the existence of spurious local minima, together with Theorem 4.12, we eliminate a large number of such points.

4.4. Absence of correlation between tractability and benign landscape.

So far, we have witnessed that the instances of (FSMF) satisfying the assumptions of Theorem 3.8 are not only efficiently solvable using Algorithm 3.3: they also have a landscape with no spurious local valleys and favorable in terms of spurious local minima Theorem 4.19. The question of interest is: Is there a link between such benign landscape and the tractability of the problem? Even if the natural answer could intuitively seem to be positive, as it is the case for the full support case, we prove that this conjecture is not true. We first provide a counter example showing that tractability does not imply a benign landscape. Then, in Section 5 we provide numerical illustration of the fact that even with a benign landscape the convergence of the gradient descent method may not be straightforward.

First, we provide a sufficient condition for the *existence* of a spurious local valley in (FSMF).

THEOREM 4.20. *Consider function $L(X, Y) = \|A - XY^\top\|^2$ in (FSMF). Given two support constraints $I \in \{0, 1\}^{m \times r}$, $J \in \{0, 1\}^{n \times r}$, if there exist $i_1 \neq i_2 \in \llbracket m \rrbracket$, $j_1 \neq j_2 \in \llbracket n \rrbracket$ and $k \in \llbracket r \rrbracket$ such that (i_1, j_1) belongs to at least 2 rank-one supports, one of which is \mathcal{S}_k , and if $(i_1, j_2), (i_2, j_1), (i_2, j_2)$ belong only to \mathcal{S}_k , then:*

- 1) *There exists A such that: $L(X, Y)$ has a spurious local valley.*
- 2) *There exists A such that: $L(X, Y)$ has a spurious local minimum.*

In both cases, A can be chosen such that $A_{i_2, j_2} \neq 0$.

Remark 4.21. The property $A_{i_2, j_2} \neq 0$ is important, as it allows to build a counter-example to the mentioned conjecture, cf. Example 4.22 below.

Proof. Let $l \neq k$ be another rank-one contribution support \mathcal{S}_l that contains (i_1, j_1) . Without loss of generality, we can assume $i_1 = j_1 = 1, i_2 = j_2 = 2$ and $k = 1, l = 2$. In particular, let $I' = J' := \{(1, 1), (1, 2), (2, 1)\}$, then $I' \subseteq I, J' \subseteq J$.

- 1) We define the matrix A by block matrices as:

$$A = \begin{pmatrix} A' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \text{ where } A' = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \in \mathbb{R}^{2 \times 2}.$$

Thus, $A_{i_2, j_2} = A_{2,2} \neq 0$. The minimum of $L(X, Y) := \|A - XY^\top\|^2$ over feasible pairs is zero and it is attained at $X = \begin{bmatrix} X' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, Y = \begin{bmatrix} Y' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ where $X' = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}, Y' = \begin{bmatrix} 1 & -1 \\ 1 & 0 \end{bmatrix}$. (X, Y) is feasible since $\text{supp}(X) = \text{supp}(X') = I' \subseteq I, \text{supp}(Y) =$

$\text{supp}(Y') = J' \subseteq J$. Moreover,

$$(4.3) \quad XY^\top = \begin{pmatrix} X'Y'^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} A' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = A.$$

Using [Lemma 4.8](#) we now prove that this matrix A produces spurious local valleys for $L(X, Y)$ with any support constraints (I, J) satisfying the assumptions. In fact, since $(1, 2), (2, 1), (2, 2)$ are only in \mathcal{S}_1 and in no other support $\mathcal{S}_\ell, \ell \neq 1$, one can easily check that for every feasible pair (X, Y) we have:

$$(4.4) \quad (XY^\top)_{i,j} = X_{i,1}Y_{j,1}, \quad \forall (i, j) \in \{(1, 2), (2, 1), (2, 2)\}.$$

Thus, every feasible pair (X^*, Y^*) reaching the global optimum $\|A - X^*(Y^*)^\top\| = 0$ must satisfy $X_{2,1}^*Y_{1,1}^* = X_{1,1}^*Y_{2,1}^* = X_{2,1}^*Y_{2,1}^* = 1$. This implies $X_{1,1}^*Y_{1,1}^* = (X_{2,1}^*Y_{1,1}^*)(X_{1,1}^*Y_{2,1}^*)/(X_{2,1}^*Y_{2,1}^*) = 1$. Moreover, such an optimum feasible pair must also satisfy $A_{11} = (X^*(Y^*)^\top)_{1,1} = \sum_p X_{1,p}^*Y_{1,p}^*$, hence $\sum_{p \neq 1} X_{1,p}^*Y_{1,p}^* = A_{1,1} - X_{1,1}^*Y_{1,1}^* = -1$.

To show the existence of a spurious local valley we use [Lemma 4.8](#) and consider the set $\tilde{S}_\sigma = \{(X, Y) \mid \text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J, \sum_{p \neq 1} X_{1,p}Y_{1,p} = \sigma\}$. We will show that $S_1 := \tilde{S}_{-1}, S_2 := \tilde{S}_1, S_3 := \tilde{S}_5$ satisfy the assumptions of [Lemma 4.8](#).

To compute $\inf_{(X,Y) \in \tilde{S}_i} L(X, Y)$, we study $g(\sigma) := \inf_{(X,Y) \in \tilde{S}_\sigma} L(X, Y)$. Denoting $Z = \begin{bmatrix} \mathbf{I}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \{0, 1\}^{m \times n}$ we have:

$$\begin{aligned} g(\sigma) &= \inf_{(X,Y) \in \tilde{S}_\sigma} \|A - XY^\top\|^2 \\ &\geq \inf_{(X,Y) \in \tilde{S}_\sigma} \|(A - XY^\top) \odot Z\|^2 \\ &\stackrel{(4.4)}{=} \inf_{(X,Y) \in \tilde{S}_\sigma} \left\| \begin{pmatrix} A_{1,1} - \sigma - X_{1,1}Y_{1,1} & A_{1,2} - X_{1,1}Y_{2,1} \\ A_{2,1} - X_{2,1}Y_{1,1} & A_{2,2} - X_{2,1}Y_{2,1} \end{pmatrix} \right\|^2 \\ &= \inf_{X_{1,1}, X_{2,1}, Y_{1,1}, Y_{2,1}} \left\| \begin{pmatrix} -\sigma - X_{1,1}Y_{1,1} & 1 - X_{1,1}Y_{2,1} \\ 1 - X_{2,1}Y_{1,1} & 1 - X_{2,1}Y_{2,1} \end{pmatrix} \right\|^2 \end{aligned}$$

Besides [Equation \(4.4\)](#), the third equality exploits the fact that $(XY^\top)_{1,1} = \sum_p X_{1,p}Y_{1,p} = X_{1,1}Y_{1,1} + \sigma$. The last quantity is the loss of the best rank-one approximation of $\tilde{A} = \begin{bmatrix} -\sigma & 1 \\ 1 & 1 \end{bmatrix} \in \mathbb{R}^{2 \times 2}$. Since this is a 2×2 symmetric matrix, its eigenvalues can be computed as the solutions of a second degree polynomial, leading to an analytic expression of this last quantity as: $\frac{2(\sigma+1)^2}{(\sigma^2+3) + \sqrt{(\sigma^2+3)^2 - 4(\sigma+1)^2}}$.

Moreover, this infimum can be attained if $[X_{1,1}, X_{2,1}] = [Y_{1,1}, Y_{2,1}]$ is the first eigenvector of \tilde{A} and the other coefficients of X, Y are set to zero. Therefore,

$$g(\sigma) = \frac{2(\sigma+1)^2}{(\sigma^2+3) + \sqrt{(\sigma^2+3)^2 - 4(\sigma+1)^2}}.$$

We can now verify that S_1, S_2, S_3 satisfy all the conditions of [Lemma 4.8](#).

- 1) The minimum value of L is zero. As shown above, it is only attained with $\sum_{p \neq 1} X_{1,p}^*Y_{1,p}^* = -1$ as shown. Thus, the global minima belong to $S_1 = \tilde{S}_{-1}$.
- 2) For any feasible path $r : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : t \rightarrow (X(t), Y(t))$ we have $\sigma_r(t) = \sum_{p \neq 1} X(t)_{1,p}Y(t)_{1,p}$ is also continuous. If $(X(0), Y(0)) \in S_3 = \tilde{S}_5$ and $(X(1), Y(1)) \in S_1 = \tilde{S}_{-1}$ then $\sigma_r(0) = 5$ and $\sigma_r(1) = -1$, hence by the

Mean Value Theorem, there must exist $t \in (0, 1)$ such that $\sigma_r(t) = 1$, which means $(X(t), Y(t)) \in S_2 = \tilde{S}_1$.

3) Since one can check numerically that $g(1) > g(5) > g(-1)$, we have

$$\inf_{(X,Y) \in S_2} L(X, Y) > \inf_{(X,Y) \in S_3} L(X, Y) > \inf_{(X,Y) \in S_1} L(X, Y).$$

The proof is concluded with the application of [Lemma 4.8](#).

2) We define the matrix A by block matrices as:

$$A = \begin{pmatrix} A' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \text{ where } A' = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \in \mathbb{R}^{2 \times 2}.$$

where $a > b > 0$. Thus, $A_{i_2, j_2} = A_{2,2} \neq 0$. It is again evident that $\inf_{X,Y} \|A - XY^\top\|^2 = 0$ (taking $X = \begin{bmatrix} X' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, $Y = \begin{bmatrix} Y' & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$ where $X' = \begin{bmatrix} 0 & a \\ b & 0 \end{bmatrix}$, $Y' = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ and with the same proof as in [Equation \(4.3\)](#), we have $XY^\top = A$).

Now, we will consider $\tilde{X} = \begin{bmatrix} a & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, $\tilde{Y} = \begin{bmatrix} 1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$. Since $L(\tilde{X}, \tilde{Y}) = b^2 > 0$ it cannot be a global minimum. We will show that (\tilde{X}, \tilde{Y}) is indeed a local minimum, which will thus imply that (\tilde{X}, \tilde{Y}) is a spurious local minimum. For each feasible pair (X, Y) we have:

$$\begin{aligned} \|A - XY^\top\|^2 &= \sum_{i,j} (A_{i,j} - (XY^\top)_{i,j})^2 \\ &\geq (A_{2,1} - (XY^\top)_{2,1})^2 + (A_{1,2} - (XY^\top)_{1,2})^2 + (A_{2,2} - (XY^\top)_{2,2})^2 \\ &\stackrel{(4.4)}{=} (X_{2,1}Y_{1,1})^2 + (X_{1,1}Y_{2,1})^2 + (b - X_{2,1}Y_{2,1})^2 \\ &\geq 2(X_{1,1}Y_{1,1})|X_{2,1}Y_{2,1}| + (X_{2,1}Y_{2,1})^2 - 2bX_{2,1}Y_{2,1} + b^2 \\ &\geq 2(X_{1,1}Y_{1,1} - b)|X_{2,1}Y_{2,1}| + b^2. \end{aligned}$$

where in the third line we used that for $u = |X_{2,1}Y_{1,1}|$, $v = |X_{1,1}Y_{2,1}|$, since $(u - v)^2 \geq 0$ we have $u^2 + v^2 \geq 2uv$. Since $\tilde{X}_{1,1}\tilde{Y}_{1,1} = a > b$, there exists a neighborhood of (\tilde{X}, \tilde{Y}) such that $X_{1,1}Y_{1,1} - b > 0$ for all (X, Y) in that neighbourhood. Since $|X_{2,1}Y_{2,1}| \geq 0$ in this neighborhood it follows that $\|A - XY^\top\|^2 \geq b^2 = L(\tilde{X}, \tilde{Y}) > 0$ in that neighborhood. This concludes the proof. \square

We can now exhibit the announced counter-example to the mentioned conjecture:

Example 4.22. Consider an instance of (FSMF) with $I = J = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. This pair (I, J) satisfies the assumptions of [Theorem 4.20](#) with $i_1 = 1, i_2 = 2, j_1 = 1, j_2 = 2$. Thus, with well chosen $A \in \mathbb{R}^{2 \times 2}$, $A = (A_{i,j}), 1 \leq i, j \leq 2$ such that $A_{2,2} \neq 0$, the landscape admits spurious objects. On the other hand, the problem is tractable for every $A \in \mathbb{R}^{2 \times 2}$ with $A_{2,2} \neq 0$. Indeed, $\inf_{\text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J} L(X, Y) = 0$ with optimal factors analytically given by: $X = \begin{bmatrix} 1 & A_{1,2}/A_{2,2} \\ 0 & 1 \end{bmatrix}$, $Y = \begin{bmatrix} A_{1,1} - A_{1,2}A_{2,1}/A_{2,2} & A_{2,1} \\ 0 & A_{2,2} \end{bmatrix}$. When $A_{2,2} = 0$, the infimum of $L(X, Y)$ might not be achievable, see [Remark A.1](#).

The existence of spurious local valleys shown in [Theorem 4.20](#) highlights the importance of initialization: if an initial point is already inside a spurious valley, first-order methods cannot escape this suboptimal area. An optimist may wonder if there nevertheless exist a smart initialization that avoids all spurious local valleys initially. The answer is positive, as shown in the following theorem.

THEOREM 4.23. *Given any I, J, A such that the infimum of (FSMF) is attained, every initialization $(X, \mathbf{0})$, $\text{supp}(X) \subseteq I$ (or symmetrically $(\mathbf{0}, Y)$, $\text{supp}(Y) \subseteq J$) is not in any spurious local valley. In particular, $(\mathbf{0}, \mathbf{0})$ is never in any spurious local valley.*

Proof. Let (X^*, Y^*) be a minimizer of (FSMF), which exists due to our assumptions. We only prove the result for the initialization $(X, \mathbf{0})$, $\text{supp}(X) \subseteq I$. The case of the initialization $(\mathbf{0}, Y)$, $\text{supp}(Y) \subseteq J$ can be dealt with similarly.

To prove the theorem, it is sufficient to construct $f(t) = (X_f(t), Y_f(t)) : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$ a feasible path such that:

- 1) $f(0) = (X, \mathbf{0})$.
- 2) $f(1) = (X^*, Y^*)$.
- 3) $L \circ f$ is non-increasing w.r.t t .

Indeed, if such f exists, the sublevel set corresponding to $L(X, \mathbf{0})$ has both $(X, \mathbf{0})$ and (X^*, Y^*) in the same path-connected components (since $L \circ f$ is non-increasing).

We will construct such a function feasible path f as a concatenation of two functions feasible paths $f_1 : [0, 1/2] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$, $f_2 : [1/2, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$, defined as follows:

- 1) $f_1(t) = ((1 - 2t)X + 2tX^*, \mathbf{0})$.
- 2) $f_2(t) = (X^*, (2t - 1)Y^*)$.

It is obvious that $f(0) = f_1(0) = (X, \mathbf{0})$ and $f(1) = f_2(1) = (X^*, Y^*)$. Moreover f is continuous since $f_1(1/2) = f_2(1/2) = (X^*, \mathbf{0})$. Also, $L \circ f$ is non-increasing on $[0, 1]$ since:

- 1) $L(f_1(t)) = \|A - ((1 - 2t)X + 2tX^*)\mathbf{0}^\top\|^2 = \|A\|^2$ is constant for $t \in [0, 1/2]$.
- 2) $L(f_2(t)) = \|A - (2t - 1)X^*Y^*\|^2$ is convex w.r.t t . Moreover, it attains a global minimum at $t = 1$ (since we assume that (X^*, Y^*) is a global minimizer of (FSMF)). As a result, $t \mapsto L(f_2(t))$ is non-increasing on $[1/2, 1]$. \square

Yet, such an initialization does not guarantee that first-order methods converge to a global minimum. Indeed, while in the proof of this result we do show that there exists a feasible path joining this “smart” initialization to an optimal solution without increasing the loss function, the value of the objective function is “flat” in the first part of this feasible path. Thus, even if such initialization is completely outside any spurious local valley, it is not clear whether local information at the initialization allow to “guide” optimization algorithms towards the global optimum to blindly find such a path. Indeed, first-order methods are not bound to follow our constructive continuous path. In the next section we further elaborate on the importance of the starting guess for local optimization methods.

5. Numerical illustration: landscape’s properties and convergence of gradient descent. As shown in Subsection 4.3, (FSMF) has a “good” landscape under the assumptions of Theorem 3.8. This might suggest that, from a random initialization (or from a “smart” one as suggested by Theorem 4.23), popular optimization methods such as gradient descent might easily be able to return the globally optimal solution. The situation is in fact more tricky. Actually, the effectiveness of those methods in this specific case has never been shown in practice. Thus, this section shows the empirical performance of gradient descent in tackling the problem of matrix factorization with fixed support.

Consider the following minimalistic instance of (FSMF): $A = \begin{bmatrix} 0 & 1 \end{bmatrix}$, $I = \begin{bmatrix} 1 \end{bmatrix}$, $J = \begin{bmatrix} 1 & 1 \end{bmatrix}$. This instance can easily be checked to satisfy the assumptions of Theorem 3.3, thus its landscape is free of spurious objects by Theorem 4.12. The infimum of this instance is zero, attained by solutions of the form $X^* = \begin{bmatrix} a \end{bmatrix}$, $Y^* = \begin{bmatrix} 0 & b \end{bmatrix}$ with $ab = 1$.

We perform gradient descent for this instance. We denote $X = \begin{bmatrix} x \end{bmatrix}$, $Y = \begin{bmatrix} y_1 & y_2 \end{bmatrix}$ and we define $g(Y) := g(y_1, y_2) = \min_X L(X, Y) = \min_x (xy_1)^2 + (1 - xy_2)^2$. Empirical experiments show that the application of gradient descent to $L(X, Y)$ is very well approximated by the application of gradient descent to $g(y_1, y_2)$. We consider then

this procedure, that allows us for instance to have a 3D visualization as in [Figure 5](#) (this is not possible for the original problem that has 3 parameters in total). [Figure 5b](#) (the loss surface of function $g(y_1, y_2)$) also shows visual proof of the fact that the landscape has no spurious local object, as proved in [Theorem 4.12](#).

With fixed y_1, y_2 , g is a simple quadratic function w.r.t x . Solving the quadratic minimization problem, we have $g(y_1, y_2) = y_1^2/(y_1^2 + y_2^2)$ and it is attained with $x = y_2/(y_1^2 + y_2^2)$. We consider two initializations where X is already the optimum given Y : $X_0 = [1/170]$, $Y_0 = [40 \ 10]$ and $X_1 = [0.2]$, $Y_1 = [2 \ 1]$, which both satisfy the condition $x = y_2/(y_1^2 + y_2^2)$. The learning rate α is chosen by backtracking line search, satisfying the Armijo condition [\[16\]](#)⁴.

From [Figure 5a](#) it is clear that the performance of the gradient descent is deeply affected by the choice of the initial guess, despite the absence of spurious objects in the landscape. Indeed, [Figure 5c](#) presents the surface of the gradient of $g(y_1, y_2)$ and shows that the sequence generated starting from (X_0, Y_0) (blue line on the right) resides completely inside an area with very small gradient. In addition, the landscape around (X_0, Y_0) is very flat ([Figure 5b](#)). Therefore, gradient descent has a lot of difficulties to converge to the optimum. In contrast, (X_1, Y_1) lies in an area with larger gradient and without any flat surrounding area. As a consequence, its corresponding sequence achieves optimality much faster. Initializing with $X'_0 = 0$ and Y_0 (resp. with $X'_1 = 0$ and Y_1) yields the same behavior.

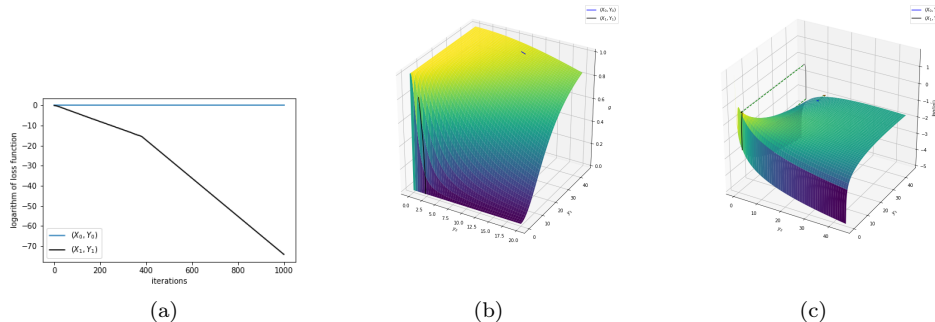


FIG. 5. (a) Evolution of the logarithm of $L(X, Y)$ with two different initializations. (b) The surface of $g(y_1, y_2) = \min_X L(X, Y)$ (c) The surface of $\log \|\nabla g(y_1, y_2)\|$. Trajectories of gradient descent from (X_0, Y_0) after 10^5 iterations (blue) and from (X_1, Y_1) after 10^3 iterations (black). Figure (c) contains the projections of the gradient norms of two trajectories.

The example shows that the effectiveness of gradient descent for (FSMF) heavily depends on initialization, which is not evident to choose. In contrast, our [Algorithm 3.1](#) does not require to tune any hyper-parameter.

6. Conclusion. In this paper, we studied the problem of two-layer matrix factorization with fixed support. We showed that this problem is NP-hard in general. Nevertheless, certain structured supports allow for an efficient solution algorithm. Furthermore, we also showed the non-existence of spurious objects in the landscape of function $L(X, Y)$ of (FSMF) with these support constraints. Although it would have

⁴For the problem $\min_x f(x)$ Armijo condition requires α to satisfy $f(x - \alpha \nabla f(x)) \leq f(x) - \alpha c \|\nabla f(x)\|^2$, we set $c = 10^{-4}$.

seemed natural to assume an equivalence between tractability and benign landscape of (FSMF), we also show a counter-example that contradicts this conjecture. That shows that there is still room for improvement of the current tools (spurious objects) to characterize the tractability of an instance. We have also shown numerically the limitations of state-of-the-art first-order optimization methods in this context. In particular the convergence is highly affected from the choice of the hyper-parameters, even in the absence of spurious objects, while the proposed method does not need to tune anything. We refer the reader to [11] where we propose an extension of Algorithm 3.1 to fixed-support multilayer sparse factorization and show the superiority of the resulting method in terms of both accuracy and speed compared to the state of the art [2].

REFERENCES

- [1] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, USA, 2004.
- [2] T. DAO, A. GU, M. EICHHORN, A. RUDRA, AND C. RE, *Learning fast algorithms for linear transforms using butterfly factorizations*, in Proceedings of the 36th International Conference on Machine Learning, vol. 97, 2019, pp. 1517–1527, <http://proceedings.mlr.press/v97/dao19a.html>.
- [3] C. ECKART AND G. YOUNG, *The approximation of one matrix by another of lower rank*, Psychometrika, 1 (1936), pp. 211–218, <https://doi.org/10.1007/BF02288367>.
- [4] S. FOUCART AND H. RAUHUT, *A Mathematical Introduction to Compressive Sensing*, Springer, 2013, <https://doi.org/10.1007/978-0-8176-4948-7>, <http://link.springer.com/10.1007/978-0-8176-4948-7>.
- [5] N. GILLIS AND F. GLINEUR, *Low-rank matrix approximation with weights or missing data is NP-hard*, SIAM Journal on Matrix Analysis and Applications, 32 (2010), <https://doi.org/10.1137/110820361>.
- [6] G. GOLUB AND W. KAHAN, *Calculating the singular values and pseudo-inverse of a matrix*, Journal of the Society for Industrial and Applied Mathematics: Series B, Numerical Analysis, 2 (1965), pp. 205–224, <http://www.jstor.org/stable/2949777>.
- [7] W. HACKBUSCH, *A sparse matrix arithmetic based on h-matrices. part i: Introduction to h-matrices*, Computing, 62 (1999), pp. 89–108.
- [8] W. HACKBUSCH AND B. N. KHOROMSKIJ, *A sparse h-matrix arithmetic. part ii: Application to multi-dimensional problems*, Computing, 64 (2000), pp. 21–47.
- [9] K. KAWAGUCHI, *Deep learning without poor local minima*, in Advances in Neural Information Processing Systems 29, 2016, pp. 586–594, <http://papers.nips.cc/paper/6112-deep-learning-without-poor-local-minima.pdf>.
- [10] N. KISHORE KUMAR AND J. SHNEIDER, *Literature survey on low rank approximation of matrices*, ArXiv preprint 1606.06511, (2016).
- [11] Q. LE, L. ZHENG, E. RICCIETTI, AND R. GRIBONVAL, *Fast learning of fast transforms, with guarantees*, tech. report, 2021.
- [12] L. LE MAGOAROU AND R. GRIBONVAL, *Chasing butterflies: In search of efficient dictionaries*, in 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2015, pp. 3287–3291, <https://doi.org/10.1109/ICASSP.2015.7178579>.
- [13] L. LE MAGOAROU AND R. GRIBONVAL, *Flexible Multi-layer Sparse Approximations of Matrices and Applications*, IEEE Journal of Selected Topics in Signal Processing, 10 (2016), pp. 688–700.
- [14] X. LI, J. LU, R. ARORA, J. HAUPT, H. LIU, Z. WANG, AND T. ZHAO, *Symmetry, saddle points, and global optimization landscape of nonconvex matrix factorization*, IEEE Transactions on Information Theory, 65 (2019), pp. 3489–3514, <https://doi.org/10.1109/TIT.2019.2898663>.
- [15] Q. NGUYEN, *On connected sublevel sets in deep learning*, in Proceedings of the 36th International Conference on Machine Learning, vol. 97, 2019, pp. 4790–4799, <http://proceedings.mlr.press/v97/nguyen19a.html>.
- [16] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, second ed., 2006.
- [17] R. PEETERS, *The maximum edge biclique problem is NP-complete*, Discrete Appl Math, 131 (2000).
- [18] R. RUBINSTEIN, A. BRUCKSTEIN, AND M. ELAD, *Dictionaries for sparse representation modeling*, Proceedings of the IEEE, 98 (2010), pp. 1045 – 1057, <https://doi.org/10.1109/>

- JPROC.2010.2040551.
- [19] R. RUBINSTEIN, M. ZIBULEVSKY, AND M. ELAD, *Double sparsity: Learning sparse dictionaries for sparse signal approximation*, IEEE Transactions on Signal Processing, 58 (2010), pp. 1553–1564, <https://doi.org/10.1109/TSP.2009.2036477>.
 - [20] Y. SHITOV, *A short proof that NMF is NP-hard*, Arxiv preprint 1605.04000, (2016).
 - [21] I. TOŠIĆ AND P. FROSSARD, *Dictionary learning*, IEEE Signal Processing Magazine, 28 (2011), pp. 27–38, <https://doi.org/10.1109/MSP.2010.939537>.
 - [22] L. N. TREFETHEN AND D. BAU III, *Numerical linear algebra*, vol. 50, SIAM, 1997.
 - [23] S. A. VAVASIS, *On the complexity of nonnegative matrix factorization*, SIAM Journal on Optimization, 20 (2010), pp. 1364–1377, <https://doi.org/10.1137/070709967>, <https://doi.org/10.1137/070709967>, <https://arxiv.org/abs/https://doi.org/10.1137/070709967>.
 - [24] L. VENTURI, A. S. BANDEIRA, AND J. BRUNA, *Spurious valleys in one-hidden-layer neural network optimization landscapes*, Journal of Machine Learning Research, 20 (2019), pp. 1–34, <http://jmlr.org/papers/v20/18-674.html>.
 - [25] Z. ZHU, D. SOUDRY, Y. C. ELДАР, AND M. WAKIN, *The global optimization geometry of shallow linear neural networks*, Journal of Mathematical Imaging and Vision, 62 (2019), pp. 279–292.

Appendix A. Proof of Lemma 2.3. Up to a transposition, we can assume WLOG that $m \geq n$. We will show that with $r = n + 1 = \min(m, n) + 1$, we can find two supports I and J satisfying the conclusion of Lemma 2.3.

To create an instance of (FSMF) (i.e., two supports I, J) that is *equivalent* to (MCPO), we define $I \in \{0, 1\}^{m \times (n+1)}$ and $J \in \{0, 1\}^{n \times (n+1)}$ as follows:

$$(A.1) \quad I_{i,j} = \begin{cases} 1 - W_{i,j} & \text{if } j \neq n \\ 1 & \text{if } j = n + 1 \end{cases}, \quad J_{i,j} = \begin{cases} 1 & \text{if } j = i \text{ or } j = n + 1 \\ 0 & \text{otherwise} \end{cases}$$

Figure 6 illustrates an example of support constraints built from W .

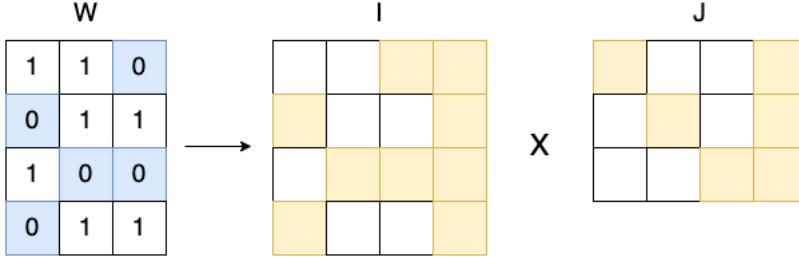


FIG. 6. Factor supports I and J constructed from the weighted matrix $W \in \{0, 1\}^{4 \times 3}$. Colored squares in I and J are positions in the supports.

We consider the (FSMF) with the same matrix A and I, J defined as in Equation (A.1). This construction (of I and J) can clearly be made in polynomial time. Consider the coefficients $(XY^\top)_{i,j}$:

- 1) If $W_{i,j} = 0$: $(XY^\top)_{i,j} = \sum_{k=1}^{n+1} X_{i,k}Y_{j,k} = X_{i,j}Y_{j,j} + X_{i,n+1}Y_{j,n+1}$ (except for $k = n + 1$, only $Y_{j,j}$ can be different from zero due to our choice of J).
- 2) If $W_{i,j} = 1$: $(XY^\top)_{i,j} = \sum_{k=1}^{n+1} X_{i,k}Y_{j,k} = X_{i,n+1}Y_{j,n+1}$ (same reason as in the previous case, in addition to the fact that $I_{i,j} = 1 - W_{i,j} = 0$).

Therefore, the following equation holds:

$$(A.2) \quad (XY^\top) \odot W = (X_{\bullet, n+1}Y_{n+1, \bullet}^\top) \odot W$$

We will prove that (FSMF) and (MCPO) share the same infimum⁵. Let $\mu_1 = \inf_{x \in \mathbb{R}^m, y \in \mathbb{R}^n} \|A - xy^\top\|_W^2$ and $\mu_2 = \inf_{\text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J} \|A - XY^\top\|^2$. It is clear that $\mu_i \geq 0 > -\infty, i = 1, 2$. Our objective is to prove $\mu_1 \leq \mu_2$ and $\mu_2 \leq \mu_1$.

- 1) Proof of $\mu_1 \leq \mu_2$: By definition of an infimum, for all $\mu > \mu_1$, there exist x, y such that $\|A - xy^\top\|_W^2 \leq \mu$. We can choose X and Y (with $\text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$) as follows: we take the last columns of X and Y equal to x and y ($X_{\bullet, n+1} = x, Y_{\bullet, n+1} = y$). For the *remaining* columns of X and Y , we choose:

$$\begin{aligned} X_{i,j} &= A_{i,j} - x_i y_j & \text{if } I_{i,j} = 1, j \leq n \\ Y_{i,j} &= 1 & \text{if } J_{i,j} = 1, j \leq n \end{aligned}$$

This choice of X and Y will make $\|A - XY^\top\|^2 = \|A - xy^\top\|_W^2 \leq \mu$. Indeed, for all (i, j) such that $W_{i,j} = 0$, we have:

$$(A - XY^\top)_{i,j} = A_{i,j} - X_{i,j}Y_{j,j} - X_{i,n+1}Y_{j,n+1} = A_{i,j} - A_{i,j} + x_i y_j - x_i y_j = 0$$

⁵We focus on the infimum instead of minimum since there are cases where the infimum is not attained, as shown in Remark A.1

Therefore, it is clear that: $(A - XY^\top) \odot (\mathbf{1} - W) = \mathbf{0}$.

$$\begin{aligned} \|A - XY^\top\|^2 &= \|(A - XY^\top) \odot W\|^2 + \|(A - XY^\top) \odot (\mathbf{1} - W)\|^2 \\ &= \|(A - XY^\top) \odot W\|^2 \\ &\stackrel{\text{(A.2)}}{=} \|(A - X_{\bullet, n+1} Y_{\bullet, n+1}^\top) \odot W\|^2 \\ &= \|(A - xy^\top) \odot W\|^2 \\ &= \|A - xy^\top\|_W^2 \end{aligned}$$

Therefore, $\mu_2 \leq \mu_1$.

- 2) Proof of $\mu_1 \leq \mu_2$: Inversely, for all $\mu > \mu_2$, there exists X, Y satisfying $\text{supp}(X) \subseteq I$, $\text{supp}(Y) \subseteq J$ such that $\|A - XY^\top\|^2 \leq \mu$. We choose $x = X_{\bullet, n+1}$, $y = Y_{\bullet, n+1}$. It is immediate that:

$$\begin{aligned} \|A - xy^\top\|_W^2 &= \|(A - xy^\top) \odot W\|^2 \\ &= \|(A - X_{\bullet, n+1} Y_{\bullet, n+1}^\top) \odot W\|^2 \\ &\stackrel{\text{(A.2)}}{=} \|(A - XY^\top) \odot W\|^2 \\ &\leq \|(A - XY^\top) \odot W\|^2 + \|(A - XY^\top) \odot (\mathbf{1} - W)\|^2 \\ &= \|A - XY^\top\|^2 \end{aligned}$$

Thus, $\|A - xy^\top\|_W^2 \leq \|A - XY^\top\|^2 \leq \mu$. We have $\mu_1 \leq \mu_2$.

This shows that $\mu_1 = \mu_2$. Moreover, the proofs of $\mu_1 \leq \mu_2$ and $\mu_2 \leq \mu_1$ also show the procedures to obtain an optimal solution of one problem with a given accuracy ϵ provided that we know an optimal solution of the other with the same accuracy.

Remark A.1. In the proof of [Lemma 2.3](#), we focus on the infimum instead of minimum since there are cases where the infimum is not attained. Indeed, consider the following instance of (FSMF) with: $A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$, $I = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, $J = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. The infimum of this problem is zero, which can be shown by choosing: $X_k = \begin{bmatrix} -k & k \\ 0 & \frac{1}{k} \end{bmatrix}$, $Y_k = \begin{bmatrix} k & k \\ 0 & \frac{1}{k} \end{bmatrix}$. In the limit, when k goes to infinity, we have:

$$\lim_{k \rightarrow \infty} \|A - X_k Y_k^\top\|^2 = \lim_{k \rightarrow \infty} \frac{1}{k^2} = 0.$$

Yet, there does not exist any couple (X, Y) such that $\|A - XY\|^2 = 0$. Indeed, any such couple would need to satisfy: $X_{1,2} Y_{2,2} = 1$, $X_{2,2} Y_{1,2} = 1$, $X_{2,2} Y_{2,2} = 0$. However, the third equation implies that either $X_{2,2} = 0$ or $Y_{2,2} = 0$, which makes either $X_{2,2} Y_{1,2} = 0$ or $X_{1,2} Y_{2,2} = 0$. This leads to a contradiction.

In fact, I and J are constructed from the weight binary matrix $W = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ (the construction is similar to one in the proof of [Lemma 2.3](#)). Problem (MCPO) with (A, W) has unattainable infimum as well.

Appendix B. Proofs for [section 3](#).

B.1. Proof of [Lemma 3.5](#). Denote $\mathcal{P}_k = \{P_1, \dots, P_k\}$, $\mathcal{S}_{\mathcal{P}_k} = \cup_{1 \leq i \leq k} \mathcal{S}_{P_i}$ for $1 \leq k \leq \ell$ and $\mathcal{S}_{P_0} = \emptyset$. In [Algorithm 3.2](#), we only assign value for $X_{R_P, P}$ and $Y_{C_P, P}$ for $P \in \mathcal{P}^*$, thus $\text{supp}(X) \subseteq I_T$, $\text{supp}(Y) \subseteq J_T$. To prove the correctness of [Algorithm 3.2](#), we show that:

$$(B.1) \quad X_{P_k} Y_{P_k}^\top = A \odot (\mathcal{S}_{\mathcal{P}_k} \setminus \mathcal{S}_{\mathcal{P}_{k-1}})$$

Thus, we have: $XY^\top = \sum_{P \in \mathcal{P}^*} X_P Y_P^\top = \sum_{k=1}^{\ell} A \odot (\mathcal{S}_{P_k} \setminus \mathcal{S}_{P_{k-1}}) = A \odot \mathcal{S}_\ell = A \odot \mathcal{S}_T = A$ (since we assume $\text{supp}(A) = \mathcal{S}_T$). We prove Equation (B.1) by induction. To ease the reading, in this proof, we denote C_{P_k}, R_{P_k} (Definition 3.4) by C_k, R_k respectively.

For $k = 1$, we have $\mathcal{S}_{P_1} \setminus \mathcal{S}_{P_0} = C_1 \times R_1$. By Line 8 and Line 11 of Algorithm 3.2, we have $X_{C_1, P_1} Y_{R_1, P_1}^\top = A'_{C_1, R_1} = A_{C_1, R_1}$ (since (X, Y) is initialized at $(\mathbf{0}, \mathbf{0})$). Therefore, $X_{P_1} Y_{P_1}^\top = A \odot (\mathcal{S}_{P_1} \setminus \mathcal{S}_{P_0})$.

Assume that Equation (B.1) holds for all $k \leq p-1, p > 1$. We prove its correctness with k . At the iteration k^{th} , we have: $A' = A - XY^\top = A - \sum_{l < k} X_{P_l} Y_{P_l}^\top = A - A \odot \mathcal{S}_{P_{k-1}} = A \odot \bar{\mathcal{S}}_{P_{k-1}}$. Therefore, $A'_{C_k, R_k} = A_{C_k, R_k} \odot (C_k \times R_k \setminus \mathcal{S}_{P_{k-1}}) = A_{C_k, R_k} \odot (\mathcal{S}_{P_k} \setminus \mathcal{S}_{P_k})$. Again, by Line 8 and Line 11 of Algorithm 3.2, we have $X_{C_k, P_k} Y_{R_k, P_k}^\top = A'_{C_k, R_k} = A_{C_k, R_k} \odot (\mathcal{S}_{P_k} \setminus \mathcal{S}_{P_k})$. Thus, $X_{P_k} Y_{P_k}^\top = A \odot (\mathcal{S}_{P_k} \setminus \mathcal{S}_{P_{k-1}})$. That implies Equation (B.1) is correct for all k .

B.2. Proof of Theorem 3.8. First, we decompose the factors X and Y using the taxonomy of indices from Definition 3.7.

DEFINITION B.1. Given I_T, J_T and $I_T^i, J_T^i, i = 1, 2$ as in Definition 3.7, consider (X, Y) a feasible point of (FSMF), we denote:

- 1) $X_T = X \odot I_T, X_{\bar{T}}^i = X \odot I_{\bar{T}}^i$, for $i = 1, 2$.
- 2) $Y_T = Y \odot I_T, Y_{\bar{T}}^i = Y \odot I_{\bar{T}}^i$, for $i = 1, 2$.

with \odot the Hadamard product between a matrix and a support constraint (introduced in subsection 1.1).

The following is a technical result.

LEMMA B.2. Given I, J support constraints of (FSMF), consider $T, \mathcal{S}_T, \mathcal{S}_{\mathcal{P}}$ as in Definition 3.2, $X_T, X_{\bar{T}}^i, Y_T, Y_{\bar{T}}^i$ as in Definition 3.6 and assume that for all $k \in \bar{T}$, \mathcal{S}'_k is rectangular. It holds:

- C1** $\text{supp}(X_T Y_T^\top) \subseteq \mathcal{S}_T$.
- C2** $\text{supp}(X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top) \subseteq \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T$.
- C3** $\text{supp}(X_{\bar{T}}^i (Y_{\bar{T}}^j)^\top) \subseteq \mathcal{S}_T, \forall 1 \leq i, j \leq 2, (i, j) \neq (1, 1)$.

Proof. We justify (C1)-(C3) as follow:

- **C1:** Since $X_T Y_T^\top = \sum_{i \in T} X_{\bullet, i} Y_{\bullet, i}^\top$, $\text{supp}(X_T Y_T^\top) \subseteq \cup_{i \in T} \mathcal{S}_k = \mathcal{S}_T$.
- **C2:** Consider the coefficient (i, j) of $(X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top$

$$((X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top)_{i, j} = \sum_k (X_{\bar{T}}^1)_{i, k} (Y_{\bar{T}}^1)_{j, k} = \sum_{(i, k) \in I_{\bar{T}}^1, (j, k) \in J_{\bar{T}}^1} X_{i, k} Y_{j, k}$$

By the definition of $I_{\bar{T}}^1, J_{\bar{T}}^1$, $(X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top_{i, j} \neq 0$ iff $(i, j) \in \cup_{\ell \in \bar{T}} R_\ell \times C_\ell = \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T$.

- **C3:** We prove for the case of $(X_{\bar{T}}^1) (Y_{\bar{T}}^2)^\top$. Others can be proved similarly.

$$(B.2) \quad ((X_{\bar{T}}^1) (Y_{\bar{T}}^2)^\top)_{i, j} = \sum_k (X_{\bar{T}}^1)_{i, k} (Y_{\bar{T}}^2)_{j, k} = \sum_{(i, k) \in I_{\bar{T}}^1, (j, k) \in J_{\bar{T}}^2} X_{i, k} Y_{j, k}$$

Since $\forall \ell \in \bar{T}, \mathcal{S}'_\ell$ is rectangular, $\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T = \cup_{\ell \in \bar{T}} \mathcal{S}'_\ell = \cup_{\ell \in \bar{T}} R_\ell \times C_\ell$. If $(i, j) \in \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T$, Equation (B.2) shows that $((X_{\bar{T}}^1) (Y_{\bar{T}}^2)^\top)_{i, j} = 0$ since there is no k such that $(i, k) \in I_{\bar{T}}^1, (j, k) \in J_{\bar{T}}^2$ due to the definition of $I_{\bar{T}}^1, J_{\bar{T}}^2$. Moreover, $\text{supp}((X_{\bar{T}}^1) (Y_{\bar{T}}^2)^\top) \subseteq \mathcal{S}_{\mathcal{P}}$ (since $\text{supp}(X_{\bar{T}}^1) \subseteq I, \text{supp}(Y_{\bar{T}}^2) \subseteq J$). Thus, it shows that $\text{supp}((X_{\bar{T}}^1) (Y_{\bar{T}}^2)^\top) \subseteq \mathcal{S}_{\mathcal{P}} \setminus (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T) = \mathcal{S}_T$. \square

Here, we present the proof of Theorem 3.8.

Proof of Theorem 3.8. Given X, Y feasible point of the input (A, I, J) , consider $X_T, Y_T, X_{\bar{T}}^i, Y_{\bar{T}}^i, i = 1, 2$ defined as in Definition B.1. Let μ_1 and μ_2 be the infimum value of (FSMF) with (A, I, J) and with $(A', I_{\bar{T}}^1, J_{\bar{T}}^1)$ ($A' = A \odot \bar{\mathcal{S}}_T$) respectively.

First, we remark that $I_{\bar{T}}^1$ and $J_{\bar{T}}^1$ satisfy the assumptions of Theorem 3.3. Indeed, it holds $\mathcal{S}_k(I_{\bar{T}}^1, J_{\bar{T}}^1) = \mathcal{S}_k(I, J) \setminus \mathcal{S}_T = \mathcal{S}'_k$ by construction. For any two indices $k, l \in \bar{T}$, the representative rank-one supports are either equal ($\mathcal{S}'_k = \mathcal{S}'_l$) or disjoint ($\mathcal{S}'_k \cap \mathcal{S}'_l = \emptyset$) by assumption. That shows why $I_{\bar{T}}^1$ and $J_{\bar{T}}^1$ satisfy the assumptions of Theorem 3.3.

Next, we prove that $\mu_1 = \mu_2$. Since $(\mathcal{S}_T, \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T, \bar{\mathcal{S}}_{\mathcal{P}})$ form a partition of $\llbracket m \rrbracket \times \llbracket n \rrbracket$, we have $C \odot D = \mathbf{0}$, $C \neq D$, $C, D \in \{\mathcal{S}_T, \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T, \bar{\mathcal{S}}_{\mathcal{P}}\}$. From the definition of A' it holds $A' \odot \bar{\mathcal{S}}_{\mathcal{P}} = A \odot \bar{\mathcal{S}}_{\mathcal{P}}$ and $A' \odot \mathcal{S}_T = \mathbf{0}$. Moreover, it holds $(X_{\bar{T}}^1)(Y_{\bar{T}}^1)^\top \odot \mathcal{S}_T \cup \bar{\mathcal{S}}_{\mathcal{P}} = \mathbf{0}$ due to C2.

Since $\text{supp}(X_T) \subseteq I_T$, $\text{supp}(X_{\bar{T}}^i) \subseteq I_{\bar{T}}$, $\text{supp}(Y_T) \subseteq J_T$, $\text{supp}(Y_{\bar{T}}^i) \subseteq J_{\bar{T}}, i = 1, 2$, the product XY^\top can be decomposed as:

$$(B.3) \quad XY^\top = X_T Y_T^\top + \sum_{1 \leq i, j \leq 2} (X_{\bar{T}}^i)(Y_{\bar{T}}^j)^\top.$$

Consider the loss function of (FSMF) with input $(A', I_{\bar{T}}^1, J_{\bar{T}}^1)$ and solution $(X_{\bar{T}}^1, Y_{\bar{T}}^1)$:

$$(B.4) \quad \begin{aligned} & \|A' - X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top\|^2 \\ &= \|(A' - X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top) \odot \mathcal{S}_T\|^2 + \|(A' - X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top) \odot (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T)\|^2 \\ &\quad + \|(A' - X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top) \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &\stackrel{\text{C2}}{=} \|(A' - (X_{\bar{T}}^1)(Y_{\bar{T}}^1)^\top) \odot \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T\|^2 + \|A' \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &\stackrel{\text{C1+C3}}{=} \|(A - X_T Y_T^\top - \sum_{1 \leq i, j \leq 2} (X_{\bar{T}}^i)(Y_{\bar{T}}^j)^\top) \odot (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T)\|^2 + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &\stackrel{\text{B.3}}{=} \|(A - XY^\top) \odot (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T)\|^2 + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \end{aligned}$$

Perform the same calculation with (A, I, J) and solution (X, Y) :

$$(B.5) \quad \begin{aligned} & \|(A - XY^\top)\|^2 \\ &= \|(A - XY^\top) \odot \mathcal{S}_T\|^2 + \|(A - XY^\top) \odot (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T)\|^2 + \|(A - XY^\top) \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &= \|(A - XY^\top) \odot \mathcal{S}_T\|^2 + \|(A - XY^\top) \odot (\mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T)\|^2 + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \end{aligned}$$

where the last equality holds since $\text{supp}(XY^\top) \subseteq \mathcal{S}_{\mathcal{P}}$. Therefore, for any feasible point (X, Y) of instance (A, I, J) , we can choose $\tilde{X} = X_{\bar{T}}^1, \tilde{Y} = Y_{\bar{T}}^1$ feasible point of $(A', I_{\bar{T}}^1, J_{\bar{T}}^1)$ such that $\|A - XY^\top\| \geq \|A' - \tilde{X}\tilde{Y}^\top\|$ (Equation (B.4) and Equation (B.5)). This shows $\mu_1 \geq \mu_2$.

On the other hand, given any feasible point (\tilde{X}, \tilde{Y}) of instance $(A', I_{\bar{T}}^1, J_{\bar{T}}^1)$, we can construct a feasible point (X, Y) for instance (A, I, J) such that $\|A - XY^\top\|^2 = \|A' - \tilde{X}'\tilde{Y}'^\top\|^2$. We construct $(X, Y) = (X_T + X_{\bar{T}}^1 + X_{\bar{T}}^2, Y_T + Y_{\bar{T}}^1 + Y_{\bar{T}}^2)$ where:

- 1) $X_{\bar{T}}^1 = \tilde{X}, Y_{\bar{T}}^1 = \tilde{Y}$,
- 2) $X_{\bar{T}}^2, Y_{\bar{T}}^2$ can be chosen arbitrarily such that $\text{supp}(X_{\bar{T}}^2) \subseteq I_{\bar{T}}^2, \text{supp}(Y_{\bar{T}}^2) \subseteq J_{\bar{T}}^2$
- 3) X_T and Y_T such that $\text{supp}(X_T) \subseteq I_T, \text{supp}(Y_T) \subseteq J_T$ and:

$$X_T Y_T^\top = (A - (X_{\bar{T}}^1 + X_{\bar{T}}^2)(Y_{\bar{T}}^1 + Y_{\bar{T}}^2)^\top) \odot \mathcal{S}_T$$

(X_T, Y_T) exists due to [Lemma 3.5](#). By [Lemma B.2](#), with this choice we have:

$$(B.6) \quad \begin{aligned} (A - XY^\top) \odot \mathcal{S}_T &\stackrel{(B.3)}{=} (A - (X_{\bar{T}}^1 + X_{\bar{T}}^2)(Y_{\bar{T}}^1 + Y_{\bar{T}}^2)^\top - X_T Y_T^\top) \odot \mathcal{S}_T \\ &\stackrel{C1}{=} (A - (X_{\bar{T}}^1 + X_{\bar{T}}^2)(Y_{\bar{T}}^1 + Y_{\bar{T}}^2)^\top) \odot \mathcal{S}_T - X_T Y_T^\top = \mathbf{0} \end{aligned}$$

Therefore $\|A - XY^\top\|^2 = \|A' - \tilde{X}\tilde{Y}^\top\|^2$ ([Equation \(B.4\)](#) and [Equation \(B.5\)](#)). Thus, $\mu_2 \geq \mu_1$. We obtain $\mu_1 = \mu_2$. In addition, given (X, Y) an optimal solution of (FSMF) with instance (A, I, J) , we have shown how to construct an optimal solution (\tilde{X}, \tilde{Y}) with instance $(A \odot \mathcal{S}_T, I_{\bar{T}}^1, J_{\bar{T}}^1)$ and vice versa. That completes our proof. \square

The following Corollary is a direct consequence of the proof of [Theorem 3.8](#).

COROLLARY B.3. *With the same assumptions and notations as in [Theorem 3.8](#), a feasible point (X, Y) (i.e., such that $\text{supp}(X) \subseteq I, \text{supp}(Y) \subseteq J$) is an optimal solution of (FSMF) if and only if:*

- 1) $(X \odot I_{\bar{T}}^1, Y \odot J_{\bar{T}}^1)$ is an optimal solution of (FSMF) with $(A \odot \bar{\mathcal{S}}_T, I_{\bar{T}}^1, J_{\bar{T}}^1)$.
- 2) The following equation holds: $(A - XY^\top) \odot \mathcal{S}_T = \mathbf{0}$

Remark B.4. In the proof of [Theorem 3.8](#), one can choose $X_{\bar{T}}^2, Y_{\bar{T}}^2$ arbitrarily. If we choose $X_{\bar{T}}^2 = \mathbf{0}, Y_{\bar{T}}^2 = \mathbf{0}$, thanks to (B.6), X_T and Y_T has to satisfy:

$$X_T Y_T^\top = (A - (X_{\bar{T}}^1 + X_{\bar{T}}^2)(Y_{\bar{T}}^1 + Y_{\bar{T}}^2)^\top) \odot \mathcal{S}_T = (A - X_{\bar{T}}^1(Y_{\bar{T}}^1)^\top) \odot \mathcal{S}_T \stackrel{C2}{=} A \odot \mathcal{S}_T$$

Appendix C. Proofs for a key lemma. In this section, we will introduce an important technical lemma. It is used extensively for the proof of the tractability and the landscape of (FSMF) under the assumptions of [Theorem 3.8](#), cf. [Appendix D.4](#).

LEMMA C.1. *Consider I, J support constraints of (FSMF) such that $\mathcal{P}^* = \mathcal{P}$. For any CEC-full-rank feasible point (X, Y) and continuous function $g : [0, 1] \rightarrow \mathbb{R}^{m \times n}$ satisfying $\text{supp}(g(t)) \subseteq \mathcal{S}_T$ ([Definition 3.4](#)) and $g(0) = XY^\top$, there exists a feasible continuous function $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

- A1** $f(0) = (X_T, Y_T)$.
- A2** $g(t) = X_f(t)Y_f(t)^\top, \forall t \in [0, 1]$.
- A3** $\|f(z) - f(t)\|^2 \leq \mathcal{C}\|g(z) - g(t)\|^2, \forall t, z \in [0, 1]$.

where $\mathcal{C} = \max_{P \in \mathcal{P}^*} \left(\max \left(\left\| X_{R_P, P}^\dagger \right\|^2, \left\| Y_{C_P, P}^\dagger \right\|^2 \right) \right)$ (D^\dagger and $\|D\|$ denote the pseudo-inverse and operator norm of a matrix D respectively).

[Lemma C.1](#) consider the case where \mathcal{P} only contains CECs. Later in other proofs, we will control the factors (X, Y) by decomposing $X = X_T + X_{\bar{T}}$ (and $Y = Y_T + Y_{\bar{T}}$) (T, \bar{T} defined in [Definition 3.4](#)) and manipulate (X_T, Y_T) and $(X_{\bar{T}}, Y_{\bar{T}})$ separately. Since the supports of (X_T, Y_T) satisfy [Lemma C.1](#), it provides us a tool to work with (X_T, Y_T) .

The proof of [Lemma C.1](#) is carried out by induction. We firstly introduce and prove two other lemmas: [Lemma C.2](#) and [Lemma C.3](#). While [Lemma C.2](#) is [Lemma C.1](#) without support constraints, [Lemma C.3](#) is [Lemma C.1](#) where $|\mathcal{P}^*| = 1$.

LEMMA C.2. *Let $X \in \mathbb{R}^{m \times r}, Y \in \mathbb{R}^{n \times r}, \min(m, n) \leq r$ and assume that X or Y has full row rank. Given any continuous function $g : [0, 1] \rightarrow \mathbb{R}^{m \times n}$ in which $g(0) = XY^\top$, there exists a continuous function $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

- 1) $f(0) = (X, Y)$.
- 2) $g(t) = X_f(t)Y_f(t)^\top, \forall t \in [0, 1]$.

3) $\|f(z) - f(t)\|^2 \leq \mathcal{C}\|g(z) - g(t)\|^2, \forall t, z \in [0, 1]$.
 where $\mathcal{C} = \max\left(\|X^\dagger\|^2, \|Y^\dagger\|^2\right)$.

Proof. WLOG, we can assume that X has full row rank. We define f as:

$$(C.1) \quad \begin{aligned} X_f(t) &= X \\ Y_f(t) &= Y + (g(t) - g(0))^\top (XX^\top)^{-1}X = Y + (X^\dagger(g(t) - g(0)))^\top \end{aligned}$$

where $X^\dagger = X^\top(XX^\top)^{-1}$ the pseudo-inverse of X . The function Y_f is well-defined due to the assumption of X being full row rank. It is immediate for the first two constraints. Since $\|f(z) - f(t)\|^2 = \|Y_f(z) - Y_f(t)\|^2 = \|X^\dagger(g(z) - g(t))\|^2$, the third one is also satisfied as:

$$\|f(z) - f(t)\|^2 = \|X^\dagger(g(z) - g(t))\|^2 \leq \|X^\dagger\|^2 \|g(z) - g(t)\|^2 \leq \mathcal{C}\|g(z) - g(t)\|^2 \quad \square$$

LEMMA C.3. Consider I, J support of (FSMF) where $\mathcal{P}^\star = \mathcal{P} = \{P\}$, for any feasible CEC-full-rank point (X, Y) and continuous function $g : [0, 1] \rightarrow \mathbb{R}^{m \times n}$ satisfying $\text{supp}(g(t)) \subseteq \mathcal{S}_P$ (Definition 3.2) and $g(0) = XY^\top$, there exists a feasible continuous function $f : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:

B1 $f(0) = (X, Y)$.

B2 $g(t) = X_f(t)Y_f(t)^\top, \forall t \in [0, 1]$.

B3 $\|f(z) - f(t)\|^2 \leq \mathcal{C}\|g(z) - g(t)\|^2$.

where $\mathcal{C} = \max\left(\|X_{R_P, P}^\dagger\|^2, \|Y_{C_P, P}^\dagger\|^2\right)$.

Proof. WLOG, we assume that $P = [|P|], R_P = [|R_P|], C_P = [|C_P|]$. Furthermore, we can assume $|P| \geq |R_P|$ and $X_{R_P, P}$ is full row rank (due to the hypothesis and the fact that P is complete).

Since $\mathcal{P}^\star = \mathcal{P} = \{P\}$, a continuous feasible function $f(t)$ must have the form: $X_f(t) = [\tilde{X}_f(t) \ \mathbf{0}]$ and $Y_f(t) = [\tilde{Y}_f(t) \ \mathbf{0}]$ where $\tilde{X}_f : [0, 1] \rightarrow \mathbb{R}^{|R_P| \times |P|}, \tilde{Y}_f : [0, 1] \rightarrow \mathbb{R}^{|C_P| \times |P|}$ are continuous functions. f is fully determined by $(\tilde{X}_f(t), \tilde{Y}_f(t))$.

Moreover, if $g : [0, 1] \rightarrow \mathbb{R}^{m \times n}$ satisfying $\text{supp}(g(t)) \subseteq \mathcal{S}_T$, then g has to have the form: $g(t) = [\tilde{g}(t) \ \mathbf{0}]$ where $\tilde{g} : [0, 1] \rightarrow \mathbb{R}^{|R_P| \times |C_P|}$ is a continuous function.

Since $g(0) = XY^\top$, $\tilde{g}(0) = (X_{R_P, P})(Y_{C_P, P})^\top$. Thus, to satisfy each constraint **B1-B3**, it is sufficient to find \tilde{X}_f and \tilde{Y}_f such that:

B1: $\tilde{X}_f(0) = X_{R_P, P}, \tilde{Y}_f(0) = Y_{C_P, P}$.

B2: $\tilde{g}(t) = \tilde{X}_f(t)\tilde{Y}_f(t)^\top, \forall t \in [0, 1]$ because:

$$X_f(t)Y_f(t)^\top = \begin{pmatrix} \tilde{X}_f(t)\tilde{Y}_f(t)^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = \begin{pmatrix} \tilde{g}(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} = g(t)$$

B3: $\|X'(z) - X'(t)\|^2 + \|Y'(z) - Y'(t)\|^2 \leq \mathcal{C}\|A'(z) - A'(t)\|^2$ since $\|X'_f(z) - X'_f(t)\|^2 + \|Y'_f(z) - Y'_f(t)\|^2 = \|f(z) - f(t)\|^2$ and $\|A'(z) - A'(t)\| = \|g(z) - g(t)\|^2$. Such function exists thanks Lemma C.2 (since we assume $X_{R_P, P}$ has full rank). \square

Proof of Lemma C.1. We prove by induction on the size \mathcal{P} . By Lemma C.3 the result is true if $|\mathcal{P}| = 1$. Assume the result is true if $|\mathcal{P}| \leq p$. We consider the case where $|\mathcal{P}| = p + 1$. Let $P \in \mathcal{P}$ and partition \mathcal{P} into $\mathcal{P}' = \mathcal{P} \setminus \{P\}$ and $\{P\}$. Let $T' = \cup_{P' \in \mathcal{P}'} P' = T \setminus P$. Since $|\mathcal{P}'| = p$, we can use induction hypothesis. Define:

$$h_1(t) = (g(t) - X_P Y_P^\top) \odot \mathcal{S}_{\mathcal{P}'}, \quad h_2(t) = X_P Y_P^\top \odot \mathcal{S}_{\mathcal{P}'} + g(t) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}$$

We verify that the function $h_1(t)$ satisfying the hypotheses to use induction step: h_1 continuous, $\text{supp}(h_1(t)) \subseteq \mathcal{S}_{\mathcal{P}'}$ and finally $h_1(0) = (g(0) - X_P Y_P^\top) \odot \mathcal{S}_{\mathcal{P}'} = X_{T'} Y_{T'}^\top \odot \mathcal{S}_{\mathcal{P}'} = X_{T'} Y_{T'}^\top$. Using the induction hypothesis with \mathcal{P}' , there exists a function $f_1 : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f_1(t) = (X_f^1(t), Y_f^1(t))$ such that:

- 1) $\text{supp}(X_f^1(t)) \subseteq I_{T'}$, $\text{supp}(Y_f^1(t)) \subseteq J_{T'}$.
- 2) $f_1(0) = (X_{T'}, Y_{T'})$.
- 3) $h_1(t) = X_f^1(t) Y_f^1(t)^\top, \forall t \in [0, 1]$.
- 4) $\|f_1(z) - f_1(t)\|^2 \leq \mathcal{C}' \|h_1(z) - h_1(t)\|^2$.

where $\mathcal{C}' = \max_{\mathcal{P}' \in \mathcal{P}'}$ $\left(\max \left(\left\| \left\| X_{R_{\mathcal{P}'}, \mathcal{P}'}^\dagger \right\| \right\|^2, \left\| \left\| Y_{C_{\mathcal{P}'}, \mathcal{P}'}^\dagger \right\| \right\|^2 \right) \right)$.

On the other hand, $h_2(t)$ satisfies the assumptions of **Lemma C.3**: $h_2(t)$ is continuous and $\text{supp}(h_2(t)) = \text{supp}(X_P Y_P^\top \odot \mathcal{S}_{\mathcal{P}'} + g(t) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}) \subseteq \text{supp}(X_P Y_P^\top) \cup (\mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}) = \mathcal{S}_P$.

In addition, since $g(0) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'} = (X Y^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'} = (X_{T'} Y_{T'}^\top + X_P Y_P^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'} = X_P Y_P^\top \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}$, we have $h_2(0) = X_P Y_P^\top \odot \mathcal{S}_{\mathcal{P}'} + g(0) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'} = X_P Y_P^\top \odot (\mathcal{S}_{\mathcal{P}'} + \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}) = X_P Y_P^\top$. Invoking **Lemma C.3** with the singleton $\{P\}$, there exists a function $(X_f^2(t), Y_f^2(t))$ such that:

- 1) $\text{supp}(X_f^2(t)) \subseteq I_P$, $\text{supp}(Y_f^2(t)) \subseteq J_P$.
- 2) $f_2(0) = (X_P, Y_P)$.
- 3) $h_2(t) = X_f^2(t) Y_f^2(t)^\top, \forall t \in [0, 1]$.
- 4) $\|f_2(z) - f_2(t)\|^2 \leq \max \left(\left\| \left\| X_{R_P, P}^\dagger \right\| \right\|^2, \left\| \left\| Y_{C_P, P}^\dagger \right\| \right\|^2 \right) \|h_2(z) - h_2(t)\|^2$.

We construct the functions $f(t) = (X_f(t), Y_f(t))$ as:

$$X_f(t) = X_f^1(t) + X_f^2(t), \quad Y_f(t) = Y_f^1(t) + Y_f^2(t)$$

We verify the validity of this construction. f is clearly feasible due to the supports of $X_f^i(t), Y_f^i(t), i = 1, 2$. The remaining conditions are:

A1:

$$\begin{aligned} X_f(0) &= X_f^1(0) + X_f^2(0) = X_{T'} + X_P = X \\ Y_f(0) &= Y_f^1(0) + Y_f^2(0) = Y_{T'} + Y_P = Y \end{aligned}$$

A2:

$$\begin{aligned} X_f(t) Y_f(t)^\top &= X_f^1(t) Y_f^1(t)^\top + X_f^2(t) Y_f^2(t)^\top \\ &= h_1(t) + h_2(t) \\ &= (g(t) - X_P Y_P^\top) \odot \mathcal{S}_{\mathcal{P}'} + X_P Y_P^\top \odot \mathcal{S}_{\mathcal{P}'} + g(t) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'} \\ &= g(t) \odot (\mathcal{S}_{\mathcal{P}'} + \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}) = g(t) \end{aligned}$$

A3:

$$\begin{aligned} &\|f(z) - f(t)\|^2 \\ &= \|f_1(z) - f_1(t)\|^2 + \|f_2(z) - f_2(t)\|^2 \\ &\leq \mathcal{C}' \|h_1(z) - h_1(t)\|^2 + \max \left(\left\| \left\| X_{R_{\mathcal{P}'}, \mathcal{P}'}^\dagger \right\| \right\|^2, \left\| \left\| Y_{C_{\mathcal{P}'}, \mathcal{P}'}^\dagger \right\| \right\|^2 \right) \|h_2(z) - h_2(t)\|^2 \quad \square \\ &\leq \mathcal{C} (\|h_1(z) - h_1(t)\|^2 + \|h_2(z) - h_2(t)\|^2) \\ &= \mathcal{C} (\|(g(z) - g(t)) \odot \mathcal{S}_{\mathcal{P}'}\|^2 + \|(g(z) - g(t)) \odot \mathcal{S}_P \setminus \mathcal{S}_{\mathcal{P}'}\|^2) \\ &= \mathcal{C} \|g(z) - g(t)\|^2 \end{aligned}$$

Appendix D. Proofs for section 4.

D.1. Proof of Lemma 4.15. The proof relies on two intermediate results that we state first: Lemma D.1 and Corollary D.2. The idea of Lemma D.1 can be found in [24]. Since it is not formally proved as a lemma or theorem, we reprove it here for self-containedness. In fact, Lemma D.1 and Corollary D.2 are special cases of Lemma 4.15 with no support constraints and $\mathcal{P}^* = \mathcal{P} = \{P\}$ respectively.

LEMMA D.1. *Let $X \in \mathbb{R}^{R \times p}$, $Y \in \mathbb{R}^{C \times p}$, $\min(R, C) \leq p$. There exists a continuous function $f(t) = (X_f(t), Y_f(t))$ on $[0, 1]$ such that:*

- $f(0) = (X, Y)$.
- $XY^\top = X_f(t)(Y_f(t))^\top, \forall t \in [0, 1]$.
- $X_f(1)$ or $Y_f(1)$ has full row rank.

Proof. WLOG, we assume that $m \leq r$. If X has full row rank, then one can choose constant function $f(t) = (X, Y)$ to satisfy the conditions of the lemma. Therefore, we can focus on the case where $\text{rank}(X) = q < m$. WLOG, we can assume that the first q columns of X (X_1, \dots, X_q) are linearly independent. The remaining columns of X can be expressed as:

$$X_k = \sum_{i=1}^q \alpha_i^k X_i, \forall q < k \leq r$$

We define a matrix \tilde{Y} by their columns as follow:

$$\tilde{Y}_i = \begin{cases} Y_i + \sum_{k=q+1}^r \alpha_i^k Y_k & \text{if } i \leq q \\ 0 & \text{otherwise} \end{cases}$$

By construction, we have $XY^\top = X\tilde{Y}^\top$. We define the function $f_1 : [0, 1] \rightarrow \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r}$ as:

$$f_1(t) = (X, (1-t)Y + t\tilde{Y})$$

This function will not change the value of f since we have:

$$X((1-t)Y^\top + t\tilde{Y}^\top) = (1-t)XY^\top + tX\tilde{Y}^\top = XY^\top.$$

Let \tilde{X} be a matrix whose first q columns are identical to that of X and $\text{rank}(\tilde{X}) = m$. The second function f_2 defined as:

$$f_2(t) = ((1-t)X + t\tilde{X}, \tilde{Y})$$

also has their product unchanged (since first q columns of $(1-t)X + t\tilde{X}$ are constant and last $r-q$ rows of \tilde{Y} are zero). Moreover, $f_2(0) = (\tilde{X}, \tilde{Y})$ where \tilde{X} has full row rank. Therefore, the concatenation of two functions f_1 and f_2 (and shrink t by a factor of 2) are the desired function f . \square

COROLLARY D.2. *Consider I, J support constraints of (FSMF) with $\mathcal{P}^* = \mathcal{P} = \{P\}$. There is a feasible continuous function $f : [0, 1] \mapsto \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} : f(t) = (X_f(t), Y_f(t))$ such that:*

1. $f(0) = (X, Y)$;
2. $X_f(t)(Y_f(t))^\top = XY^\top, \forall t \in [0, 1]$;
3. $(X_f(1))_{R,P}$ or $(Y_f(1))_{C,P}$ has full row rank.

Proof of Corollary D.2. WLOG, up to permuting columns, we can assume $P = \llbracket P \rrbracket$, $R_P = \llbracket R_P \rrbracket$ and $C_P = \llbracket C_P \rrbracket$ (R_P and C_P are defined in Definition Definition 3.2). A feasible function $f = (X_f(t), Y_f(t))$ has the form:

$$X_f(t) = \begin{pmatrix} \tilde{X}_f(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, Y_f(t) = \begin{pmatrix} \tilde{Y}_f(t) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$

where $\tilde{X}_f : [0, 1] \mapsto \mathbb{R}^{R_P \times P}$, $\tilde{Y}_f : [0, 1] \mapsto \mathbb{R}^{C_P \times P}$.

Since P is a CEC, we have $p \geq \min(R_P, C_P)$. Hence we can use Lemma D.1 to build $(\tilde{X}_f(t), \tilde{Y}_f(t))$ satisfying all conditions of Lemma D.1. Such $(\tilde{X}_f(t), \tilde{Y}_f(t))$ fully determines f and make f our desirable function. \square

Proof of Lemma 4.15. First, we decompose X and Y as:

$$X = X_{\bar{T}} + \sum_{P \in \mathcal{P}^*} X_P, \quad Y = Y_{\bar{T}} + \sum_{P \in \mathcal{P}^*} Y_P$$

Since \bar{T} and $P \in \mathcal{P}^*$ form a partition of $\llbracket r \rrbracket$, the product XY^\top can be written as:

$$XY^\top = X_{\bar{T}}Y_{\bar{T}}^\top + \sum_{P \in \mathcal{P}^*} X_P Y_P^\top.$$

For each $P \in \mathcal{P}^*$, (I_P, J_P) contains one CEC. By applying Corollary D.2, we can build continuous functions $(X_f^P(t), Y_f^P(t))$, $\text{supp}(X_f^P(t)) \subseteq I_P$, $\text{supp}(Y_f^P(t)) \subseteq J_P$, $\forall t \in [0, 1]$ such that:

1. $(X_f^P(0), Y_f^P(0)) = (X_P, Y_P)$.
2. $X_f^P(t)(Y_f^P(t))^\top = X_P Y_P^\top, \forall t \in [0, 1]$.
3. $(X_f^P(1))_{R_P, P}$ or $(Y_f^P(1))_{C_P, P}$ has full row rank.

Our desirable $f(t) = (X_f(t), Y_f(t))$ is defined as:

$$X_f(t) = X_{\bar{T}} + \sum_{P \in \mathcal{P}^*} X_f^P(t), \quad Y(t) = Y_{\bar{T}} + \sum_{P \in \mathcal{P}^*} Y_f^P(t)$$

To conclude, it is immediate to check that $f = (X_f(t), Y_f(t))$ is feasible, $f(0) = (X, Y)$, $f(1)$ is CEC-full-rank and $X_f(t)Y_f(t)^\top = XY^\top, \forall t \in [0, 1]$. \square

D.2. Proof of Lemma 4.16. Denote $Z = XY^\top$, we construct f such that $X_f(t)Y_f(t)^\top = B(t)$, where $B(t) = Z \odot \bar{\mathcal{S}}_T + (At + Z(1-t)) \odot \mathcal{S}_T$. Such function f makes $L(X_f(t), Y_f(t))$ non-increasing since:

$$(D.1) \quad \begin{aligned} \|A - X_f(t)Y_f(t)^\top\|^2 &= \|A - B(t)\|^2 \\ &= \|(A - Z) \odot \bar{\mathcal{S}}_T\|^2 + (1-t)^2 \|(A - Z) \odot \mathcal{S}_T\|^2 \end{aligned}$$

Thus, the rest of the proof is devoted to show that such a function f exists by using Lemma C.1. Consider the function $g(t) = B(t) - X_{\bar{T}}(Y_{\bar{T}})^\top$. We have that $g(t)$ is continuous, $g(0) = B(0) - X_{\bar{T}}(Y_{\bar{T}})^\top = Z - X_{\bar{T}}(Y_{\bar{T}})^\top = X_T(Y_T)^\top$ and:

$$\begin{aligned} g(t) \odot \bar{\mathcal{S}}_T &= (B(t) - X_{\bar{T}}(Y_{\bar{T}})^\top) \odot \bar{\mathcal{S}}_T \\ &= (Z - X_{\bar{T}}(Y_{\bar{T}})^\top) \odot \bar{\mathcal{S}}_T \\ &= (X_T Y_T^\top) \odot \bar{\mathcal{S}}_T = \mathbf{0} \end{aligned}$$

which shows $\text{supp}(g(t)) \subseteq \mathcal{S}_T$. Since (X_T, Y_T) is CEC-full-rank (by our assumption, (X, Y) is CEC-full-rank), invoking Lemma C.1 with (I_T, J_T) , there exists $f^T(t) = (X_f^T(t), Y_f^T(t))$ such that:

D1 $\text{supp}(X_f^T(t)) \subseteq I_T, \text{supp}(Y_f^C(t)) \subseteq J_T$.

D2 $f^T(0) = (X_T, Y_T)$.

D3 $g(t) = X_f^T(t)(Y_f^T(t))^\top, \forall t \in [0, 1]$.

We can define our desired function $f(t) = (X_f(t), Y_f(t))$ as:

$$X_f(t) = X_T + X_f^T(t), \quad Y = Y_T + Y_f^T(t)$$

f is clearly feasible due to **(D1)**. The remaining condition to be checked is:

- First condition:

$$X_f(0) = X_f^T(0) + X_T = X_T + X_T = X, \quad Y_f(0) = Y_f^T(0) + Y_T = Y_T + Y_T = Y$$

- Second condition: holds thanks to [Equation \(D.1\)](#) and:

$$X_f(t)(Y_f(t))^\top = X_T Y_T^\top + X_f^C(t)(Y_f^C(t))^\top = X_T Y_T^\top + g(t) = B(t)$$

- Third condition:

$$\begin{aligned} (A - X_f(1)(Y_f(1))^\top) \odot \mathcal{S}_T &= (A - B(1)) \odot \mathcal{S}_T \\ &= (A - Z \odot \bar{\mathcal{S}}_T - A \odot \mathcal{S}_T) \odot \mathcal{S}_T = \mathbf{0} \end{aligned}$$

D.3. Proof of Lemma 4.17. Consider $X_T, X_T^i, Y_T, Y_T^i, i = 1, 2$ as in [Definition B.1](#). We redefine $A' = A \odot \bar{\mathcal{S}}_T, I' = I_T^1, J' = J_T^1$ as in [Theorem 3.8](#).

In light of [Corollary B.3](#), an optimal solution (\tilde{X}, \tilde{Y}) has the following form:

- 1) $\tilde{X}_T^1 = \tilde{X} \odot I_T^1, \tilde{Y}_T^1 = \tilde{Y} \odot J_T^1$ is an optimal solution of **(FSMF)** with (A', I', J') .
- 2) $\tilde{X}_T^2 = \tilde{X} \odot I_T^2, \tilde{Y}_T^2 = \tilde{Y} \odot J_T^2$ can be arbitrary.
- 3) $\tilde{X}_T = \tilde{X} \odot I_T, \tilde{Y}_T = \tilde{Y} \odot J_T$ satisfy:

$$\tilde{X}_T \tilde{Y}_T^\top = (A - \sum_{(i,j) \neq (1,1)} \tilde{X}_T^i \tilde{Y}_T^j)^\top \odot \mathcal{S}_T$$

Since (I', J') has its support constraints satisfying [Theorem 3.3](#) assumptions as shown in [Theorem 3.8](#), by [Theorem 4.12](#), there exists a function $(X_f^{\bar{T}}(t), Y_f^{\bar{T}}(t))$ such that:

- 1) $\text{supp}(X_f^{\bar{T}}(t)) \subseteq I_T^1, \text{supp}(Y_f^{\bar{T}}(t)) \subseteq J_T^1$.
- 2) $X_f^{\bar{T}}(0) = X_T^1, Y_f^{\bar{T}}(0) = Y_T^1$.
- 3) $L'(X_f^{\bar{T}}(t), Y_f^{\bar{T}}(t)) = \|A' - X_f^{\bar{T}}(t)Y_f^{\bar{T}}(t)^\top\|^2$ is non-increasing.
- 4) $(X_f^{\bar{T}}(1), Y_f^{\bar{T}}(1))$ is an optimal solution of the instance of **(FSMF)** with (A', I', J') .

Consider the function $g(t) = (A - (X_f^{\bar{T}}(t) + X_T^2)(Y_f^{\bar{T}}(t) + Y_T^2)^\top) \odot \mathcal{S}_T$. This construction makes $g(0) = X_T Y_T^\top$. Indeed,

$$\begin{aligned} g(0) &= (A - (X_f^{\bar{T}}(0) + X_T^2)(Y_f^{\bar{T}}(0) + Y_T^2)^\top) \odot \mathcal{S}_T \\ &= (A - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top) \odot \mathcal{S}_T \\ &\stackrel{(1)}{=} (XY^\top - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top) \odot \mathcal{S}_T \\ &\stackrel{(2)}{=} X_T Y_T^\top \end{aligned}$$

where (1) holds by the hypothesis $(A - XY^\top) \odot \mathcal{S}_T = \mathbf{0}$, and (2) holds by [Equation \(B.3\)](#) and $\text{supp}(X_T Y_T^\top) \subseteq \mathcal{S}_T$. Due to our hypothesis (X, Y) is CEC-full-rank, (X_T, Y_T) is CEC-full-rank. In addition, $g(t)$ continuous, $\text{supp}(g(t)) \subseteq \mathcal{S}_T$ and $g(0) = X_T Y_T^\top$. Invoking [Lemma C.1](#) with (I_T, J_T) , there exist functions $(X_f^C(t), Y_f^C(t))$ satisfying:

- 1) $\text{supp}(X_f^T(t)) \subseteq I_T, \text{supp}(Y_f^T(t)) \subseteq J_T$.
- 2) $f^T(0) = (X_T, Y_T)$.
- 3) $g(t) = X_f^T(t)Y_f^T(t)^\top, \forall t \in [0, 1]$.

Finally, one can define the function $X_f(t), Y_f(t)$ satisfying [Lemma 4.17](#) as:

$$X_f(t) = X_f^{\bar{T}}(t) + X_f^C(t) + X_f^2, \quad Y_f(t) = Y_f^{\bar{T}}(t) + Y_f^C(t) + Y_f^2$$

f is feasible due to the supports of $X_f^P(t), Y_f^P(t), P \in \{\bar{T}, C\}$ and X_f^2, Y_f^2 . The remaining conditions are satisfied as:

- First condition:

$$\begin{aligned} X_f(0) &= X_f^{\bar{T}}(0) + X_f^C(0) + X_f^2 = X_f^1 + X_T + X_f^2 = X \\ Y_f(0) &= Y_f^{\bar{T}}(0) + Y_f^C(0) + Y_f^2 = Y_f^1 + Y_T + Y_f^2 = Y \end{aligned}$$

- Second condition:

$$\begin{aligned} \|A - X_f(t)Y_f(t)^\top\|^2 &= \|A - X_f^T(t)(Y_f^T(t))^\top - (X_f^{\bar{T}}(t) + X_f^2)(Y_f^{\bar{T}}(t) + Y_f^2)^\top\|^2 \\ &= \|g(t) - X_f^T(t)Y_f^T(t)^\top\|^2 + \|(A - X_f^T(t)(Y_f^T(t))^\top) \odot \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &= \|(A' - X_f^{\bar{T}}(t)(Y_f^{\bar{T}}(t))^\top) \odot \mathcal{S}_{\mathcal{P}} \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_{\mathcal{P}}\|^2 \\ &\stackrel{\text{(B.4)}}{=} \|A' - X_f^{\bar{T}}(t)(Y_f^{\bar{T}}(t))^\top\|^2 \end{aligned}$$

Since $\|A' - X_f^{\bar{T}}(t)(Y_f^{\bar{T}}(t))^\top\|^2$ is non-increasing, so is $\|A - X_f(t)Y_f(t)^\top\|^2$.

- Third condition: By [Theorem 3.8](#), $(X_f(1), Y_f(1))$ is a global minimizer since $\|A - X_f(1)Y_f(1)^\top\|^2 = \|A' - X_f^{\bar{T}}(1)(Y_f^{\bar{T}}(1))^\top\|^2$ where $(X_f^{\bar{T}}(1), Y_f^{\bar{T}}(1))$ is an optimal solution of the instance of [\(FSMF\)](#) with (A', I', J') .

D.4. Proof of [Theorem 4.19](#). The following corollary is necessary for the proof of [Theorem 4.19](#).

COROLLARY D.3. *Consider I, J support constraints of [\(FSMF\)](#), such that $\mathcal{P}^* = \mathcal{P}$. Given any feasible CEC-full-rank point (X, Y) and any B satisfying $\text{supp}(B) \subseteq \mathcal{S}_{\mathcal{P}}$, there exists (\tilde{X}, \tilde{Y}) such that:*

$$\mathbf{E1} \quad \text{supp}(\tilde{X}) \subseteq I, \text{supp}(\tilde{Y}) \subseteq J$$

$$\mathbf{E2} \quad \tilde{X}\tilde{Y}^\top = B.$$

$$\mathbf{E3} \quad \|X - \tilde{X}\|^2 + \|Y - \tilde{Y}\|^2 \leq \mathcal{C}\|XY^\top - B\|^2.$$

where $\mathcal{C} = \max_{P \in \mathcal{P}^*} \left(\max \left(\left\| X_{R_P, P}^\dagger \right\|^2, \left\| Y_{C_P, P}^\dagger \right\|^2 \right) \right)$.

Proof. [Corollary D.3](#) is an application of [Lemma C.1](#). Consider the function $g(t) = (1-t)XY^\top + tB$. By construction, $g(t)$ is continuous, $g(0) = XY^\top$ and $\text{supp}(g(t)) \subseteq \text{supp}(XY^\top) \cup \text{supp}(B) = \mathcal{S}_{\mathcal{P}}$. Since (X, Y) is CEC-full-rank, there exists a feasible function $f(t) = (X_f(t), Y_f(t))$ satisfying [A1](#) - [A3](#) by using [Lemma C.1](#).

We choose $(\tilde{X}, \tilde{Y}) = (X_f(1), Y_f(1))$. The verification of constraints is as follow:

E1: f is feasible.

$$\mathbf{E2:} \quad \tilde{X}\tilde{Y}^\top = X_f(1)Y_f(1)^\top \stackrel{\mathbf{A2}}{=} g(1) = B.$$

$$\mathbf{E3:} \quad \|X - \tilde{X}\|^2 + \|Y - \tilde{Y}\|^2 \stackrel{\mathbf{A1}}{=} \|f(1) - f(0)\|^2 \stackrel{\mathbf{A3}}{\leq} \mathcal{C}\|g(0) - g(1)\|^2 \leq \mathcal{C}\|XY^\top - B\|^2. \square$$

Proof of [Theorem 4.19](#). As mentioned in the sketch of the proof, given any (X, Y) not CEC-full-rank, [Lemma 4.15](#) shows the existence of a path f along which L is constant and f connects (X, Y) to some CEC-full-rank (\tilde{X}, \tilde{Y}) . Therefore, this proof

will be entirely devoted to show that a feasible CEC-full-rank solution (X, Y) cannot be a spurious local minimum. This fact will be shown by the two following steps:

FIRST STEP: Consider the function $L(X, Y)$, we have:

$$L(X, Y) = \|A - XY^\top\|^2 = \|A - \sum_{P' \in \mathcal{P}^*} X_{P'} Y_{P'}^\top - X_{\bar{T}} Y_{\bar{T}}^\top\|^2$$

If (X, Y) is truly a local minimum, then $\forall P \in \mathcal{P}^*$, (X_P, Y_P) is also the local minimum of the following function:

$$L'(X_P, Y_P) = \|(A - \sum_{P' \neq P} X_{P'} Y_{P'}^\top - X_{\bar{T}} Y_{\bar{T}}^\top) - X_P Y_P^\top\|^2$$

where L' is equal to L but we optimize only w.r.t (X_P, Y_P) while fixing the other coefficients. In other words, (X_P, Y_P) is a local minimum of the problem:

$$\begin{aligned} & \underset{X' \in \mathbb{R}^{m \times r}, Y' \in \mathbb{R}^{n \times r}}{\text{Minimize}} && L'(X', Y') = \|B - X' Y'^\top\|^2 \\ & \text{Subject to:} && \text{supp}(X') \subseteq I_P \text{ and } \text{supp}(Y') \subseteq J_P \end{aligned}$$

where $B = A - \sum_{P' \neq P} X_{P'} Y_{P'}^\top - X_{\bar{T}} Y_{\bar{T}}^\top$. Since all columns of I_P (resp. of J_P) are identical, all rank-one contribution supports are totally overlapping. Thus, all local minima are global minima ([Theorem 4.12](#)). Global minima are attained when $X_P Y_P^\top = B \odot \mathcal{S}_P$ due to the expressivity of a CEC ([Lemma 3.5](#)). Thus, for any $P \in \mathcal{P}^*$, $\forall (i, j) \in \mathcal{S}_P$, we have:

$$0 = (B - X_P Y_P^\top)_{i,j} = (A - \sum_{P' \in \mathcal{P}^*} X_{P'} Y_{P'}^\top - X_{\bar{T}} Y_{\bar{T}}^\top)_{i,j} = (A - XY^\top)_{i,j}$$

which implies [Equation \(4.2\)](#).

SECOND STEP: In this step, we assume that [Equation \(4.2\)](#) holds. Consider $X_T, X_{\bar{T}}^i, Y_T, Y_{\bar{T}}^i, i = 1, 2$ as in [Definition 3.7](#). Let $A' = A \odot \bar{\mathcal{S}}_T, I' = I_{\bar{T}}^1, J' = J_{\bar{T}}^1$. We consider two possibilities. First, if $(X_{\bar{T}}^1, Y_{\bar{T}}^1)$ is an optimal solution of the instance of (FSMF) with (A', I', J') , by [Corollary B.3](#), (X, Y) is an optimal solution of (FSMF) with (A, I, J) (since [Equation \(4.2\)](#) holds). Hence it cannot be a spurious local minimum. We now focus on the second case, where $(X_{\bar{T}}^1, Y_{\bar{T}}^1)$ is *not* the optimal solution of the instance of (FSMF) with (A', I', J') . We show that in this case, in any neighborhood of (X, Y) , there exists a point (X', Y') such that $\text{supp}(X') \subseteq I$, $\text{supp}(Y') \subseteq J'$ and $L(X, Y) > L(X', Y')$. Thus (X, Y) cannot be a local minimum. Since $(I_{\bar{T}}^1, J_{\bar{T}}^1)$ satisfies [Theorem 3.3](#) assumptions, (FSMF) has no spurious local minima ([Theorem 4.12](#)). As $(X_{\bar{T}}^1, Y_{\bar{T}}^1)$ is not an optimal solution, it cannot be a local minimum either, i.e., in any neighborhood of $(X_{\bar{T}}^1, Y_{\bar{T}}^1)$, there exists (\tilde{X}, \tilde{Y}) with $\text{supp}(\tilde{X}_{\bar{T}}^1) \subseteq I', \text{supp}(\tilde{Y}_{\bar{T}}^1) \subseteq J'$ and

$$(D.2) \quad \|A' - X_{\bar{T}}^1 (Y_{\bar{T}}^1)^\top\|^2 > \|A' - \tilde{X}_{\bar{T}}^1 (\tilde{Y}_{\bar{T}}^1)^\top\|^2$$

By [Equation \(B.4\)](#), we have:

$$(D.3) \quad \begin{aligned} \|A' - (X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top\|^2 &= \|(A - (X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_P\|^2 \\ \|A' - (\tilde{X}_{\bar{T}}^1) (\tilde{Y}_{\bar{T}}^1)^\top\|^2 &= \|(A - (\tilde{X}_{\bar{T}}^1) (\tilde{Y}_{\bar{T}}^1)^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_P\|^2 \end{aligned}$$

By [Equation \(D.2\)](#) and [Equation \(D.3\)](#) we have:

$$(D.4) \quad \|(A - (X_{\bar{T}}^1) (Y_{\bar{T}}^1)^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 > \|(A - \tilde{X}_{\bar{T}}^1 (\tilde{Y}_{\bar{T}}^1)^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2$$

Consider the matrix: $B := \left(A - (\tilde{X}_T^1 + X_T^2)(\tilde{Y}_T^1 + Y_T^2)^\top \right) \odot \mathcal{S}_T$. Since $\text{supp}(B) \subseteq \mathcal{S}_T$ and (X_T, Y_T) is CEC-full-rank (we assume (X, Y) is CEC-full-rank), by [Corollary D.3](#), there exists $(\tilde{X}_T, \tilde{Y}_T)$ such that:

- 1) $\text{supp}(\tilde{X}_T) \subseteq I_T, \text{supp}(\tilde{Y}_T) \subseteq J_T$.
- 2) $\tilde{X}_T \tilde{Y}_T^\top = B$.
- 3) $\|X_T - \tilde{X}_T\|^2 + \|Y_T - \tilde{Y}_T\|^2 \leq \mathcal{C} \|X_T Y_T^\top - B\|^2$.

where $\mathcal{C} = \max_{P \in \mathcal{P}^*} \left(\max \left(\left\| \left\| X_{R_P, P}^\dagger \right\| \right\|^2, \left\| \left\| Y_{C_P, P}^\dagger \right\| \right\|^2 \right) \right)$. We define the point (\tilde{X}, \tilde{Y}) as:

$$\tilde{X} = \tilde{X}_T + \tilde{X}_T^1 + X_T^2, \quad \tilde{Y} = \tilde{Y}_T + \tilde{Y}_T^1 + Y_T^2$$

The point (\tilde{X}, \tilde{Y}) still satisfies [Equation \(4.2\)](#). Indeed,

$$\begin{aligned} (D.5) \quad (A - \tilde{X} \tilde{Y}^\top) \odot \mathcal{S}_T &= \left(A - \tilde{X}_T \tilde{Y}_T^\top - (\tilde{X}_T^1 + X_T^2)(\tilde{Y}_T^1 + Y_T^2)^\top \right) \odot \mathcal{S}_T \\ &= (B - \tilde{X}_T \tilde{Y}_T^\top) \odot \mathcal{S}_T = \mathbf{0}. \end{aligned}$$

It is clear that (\tilde{X}, \tilde{Y}) satisfies $\text{supp}(\tilde{X}) \subseteq I, \text{supp}(\tilde{Y}) \subseteq J$ due to the support of its components $(\tilde{X}_T, \tilde{Y}_T), (\tilde{X}_T^1, \tilde{Y}_T^1), (X_T^2, Y_T^2)$. Moreover, we have:

$$\begin{aligned} \|A - \tilde{X} \tilde{Y}^\top\|^2 &= \|(A - \tilde{X} \tilde{Y}^\top) \odot \mathcal{S}_T\|^2 + \|(A - \tilde{X} \tilde{Y}^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_P\|^2 \\ &\stackrel{(D.5)}{=} \|(A - \tilde{X}_T \tilde{Y}_T^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_P\|^2 \\ &\stackrel{(D.4)}{<} \|(A - X_T Y_T^\top) \odot \mathcal{S}_P \setminus \mathcal{S}_T\|^2 + \|A \odot \bar{\mathcal{S}}_P\|^2 \\ &= \|A - X Y^\top\|^2. \end{aligned}$$

Lastly, we show that (\tilde{X}, \tilde{Y}) can be chosen arbitrarily close to (X, Y) by choosing $(\tilde{X}_T^1, \tilde{Y}_T^1)$ close enough to (X_T^1, Y_T^1) . For this, denoting $\epsilon := \|X_T^1 - \tilde{X}_T^1\|^2 + \|Y_T^1 - \tilde{Y}_T^1\|^2$, we first compute:

$$\begin{aligned} \|X - \tilde{X}\|^2 + \|Y - \tilde{Y}\|^2 &= \|X_T - \tilde{X}_T\|^2 + \|Y_T - \tilde{Y}_T\|^2 + \|X_T^1 - \tilde{X}_T^1\|^2 + \|Y_T^1 - \tilde{Y}_T^1\|^2 \\ &\leq \mathcal{C} \|X_T Y_T^\top - B\|^2 + \epsilon \end{aligned}$$

We will bound the value $\|X_T Y_T^\top - B\|^2$. By using [Equation \(4.2\)](#), we have:

$$\begin{aligned} (A - \sum_{1 \leq i, j \leq 2} (X_T^i)(Y_T^j)^\top) \odot \mathcal{S}_T - X_T Y_T^\top &= (A - X_T Y_T^\top - \sum_{1 \leq i, j \leq 2} (X_T^i)(Y_T^j)^\top) \odot \mathcal{S}_T \\ &= (A - X Y^\top) \odot \mathcal{S}_T \stackrel{(4.2)}{=} \mathbf{0} \end{aligned}$$

Therefore, $X_T Y_T^\top = [A - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top] \odot \mathcal{S}_T$. We have:

$$\begin{aligned} \|X_T Y_T^\top - B\|^2 &= \|[A - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top] \odot \mathcal{S}_T - B\|^2 \\ &= \|[(\tilde{X}_T^1 + X_T^2)(\tilde{Y}_T^1 + Y_T^2)^\top - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top] \odot \mathcal{S}_T\|^2 \\ &\leq \|(\tilde{X}_T^1 + X_T^2)(\tilde{Y}_T^1 + Y_T^2)^\top - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top\|^2 \end{aligned}$$

When $\epsilon \rightarrow 0$, we have $\|(\tilde{X}_T^1 + X_T^2)(\tilde{Y}_T^1 + Y_T^2)^\top - (X_T^1 + X_T^2)(Y_T^1 + Y_T^2)^\top\| \rightarrow 0$. Therefore, with ϵ small enough, one have $\|X - X'\|^2 + \|Y - Y'\|^2$ can be arbitrarily small. This concludes the proof. \square

D.5. Proof for Example 4.22. Direct calculation of the Hessian of L at point (X_0, Y_0) is given by:

$$H(L)|_{(X_0, Y_0)} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 & 0 & 10 & 0 \\ 0 & 0 & 100 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 10 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 & 1 \end{pmatrix}$$

which is indeed positive semi-definite.