

Fast Algorithms for Sparse Reduced Rank Regression : Appendix

A Summary of results

We summarize the main steps of our paper, in red for **RRR** and in cyan for **SRRR**.

Corollary 5/6 : Strong convexity on cones of f/F^λ for RRR/SRRR	\Rightarrow	Corollary 8/9 : (proximal) PL inequality	\Rightarrow	Corollary 10 : Local linear convergence with Theorem 7
--	---------------	---	---------------	--

We also summarize the different results obtained.

Results	RRR ($\lambda = 0$)		SRRR ($0 < \lambda$)	
Local minima are global minima	\checkmark Lemma 2		\times	
Algorithm	cst_st	ls	cst_st	ls
Global convergence to a critical point	\checkmark Theorem 40	\checkmark Theorem 40	(*) Theorem 43	(*) Theorem 43
Local linear convergence	\checkmark Corollary 10	\checkmark Corollary 10	$\checkmark(\lambda < \bar{\lambda})$ Corollary 10	$\checkmark(\lambda < \bar{\lambda})$ Corollary 10

- cst_st : Algorithm 1 with fixed step size $t \leq \frac{1}{L_X}$.
- ls : Algorithm 1 with line search.
- (*) : All limit points of the sequence are critical points. If these limit points are local minima and if for any $S \subset \{1, \dots, p\}$ of cardinality at least r , the matrix $X_S^T Y$ is full-rank, then Algorithm 1 converges to a local minimum (see Appendix F.2).

B Additional definitions and classical results

In Sections B.1, B.2 and B.3, we give a few definitions that are used throughout the paper. We also recall classical results in Fact 12 and Fact 15. In Section B.4, we present the limiting subdifferential and a result for subanalytic functions in Lemma 20.

B.1 Strong convexity

Definition 11. Given $d > 0$, $\mu > 0$ and a convex set $\mathcal{V} \subset \mathbb{R}^d$, a function $f : x \in \mathcal{V} \mapsto f(x)$ is μ -strongly convex if :

$$\text{for all } x, y \in \mathcal{V}, t \in [0, 1], \quad f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{\mu}{2}t(1-t)\|y-x\|^2.$$

Fact 12. Given $d > 0$, $\mu > 0$, a convex set $\mathcal{V} \subset \mathbb{R}^d$ and a differentiable function $f : x \in \mathcal{V} \mapsto f(x)$, f is μ -strongly convex if and only if :

$$\text{for all } x, y \in \mathcal{V}, \quad f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2}\|y-x\|^2.$$

B.2 Smoothness and Lipschitz gradients

Definition 13. Given $d > 0$, $L > 0$ and a set $\mathcal{V} \subset \mathbb{R}^d$, we say that a differentiable function $f : x \in \mathcal{V} \mapsto f(x)$ has L -Lipschitz gradients in \mathcal{V} if :

$$\text{for all } x, y \in \mathcal{V}, \quad \|\nabla f(x) - \nabla f(y)\| \leq L\|y-x\|.$$

Definition 14. Given $d > 0$, $L > 0$ and a set $\mathcal{V} \subset \mathbb{R}^d$, we say that a function $f : x \in \mathcal{V} \mapsto f(x)$ is L -smooth in \mathcal{V} if it is differentiable and such that :

$$\text{for all } x, y \in \mathcal{V}, \quad f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2}\|y-x\|^2.$$

Fact 15. If f has L -Lipschitz gradients and \mathcal{V} is convex, then f is L -smooth. If f is convex and L -smooth, then f has L -Lipschitz gradients.

B.3 Sublevel sets

Definition 16. Given a set \mathcal{X} and a function $f : x \in \mathcal{X} \mapsto f(x)$, a set $\mathcal{V} \subset \mathcal{X}$ is called a sublevel set of the function f if there is $c \in \mathbb{R}$ such that :

$$\mathcal{V} = \{x \in \mathcal{X}, f(x) \leq c\}.$$

B.4 Subdifferentials, graph continuity and the Kurdyka-Łojasiewicz property

Definition 17. Given a real-valued extended function $F : \mathbb{R}^d \mapsto \mathbb{R} \cup \{\infty\}$, let

$$\text{dom } F := \{x \in \mathbb{R}^d \mid F(x) < \infty\}$$

denote its domain. For each $x \in \text{dom } F$, the Fréchet subdifferential of F at x , written $\hat{\partial}F(x)$, is the set of vectors $v \in \mathbb{R}^d$ which satisfy

$$\liminf_{y \neq x, y \rightarrow x} \frac{1}{\|y-x\|} [F(y) - F(x) - \langle v, y-x \rangle] \geq 0.$$

When $x \notin \text{dom } F$, we set $\hat{\partial}F(x) = \emptyset$. Given $x \in \mathbb{R}^d$, The limiting-subdifferential $\partial F(x)$ is defined as

$$\partial F(x) := \left\{ v \in \mathbb{R}^d \mid \exists x^k \rightarrow x, f(x^k) \rightarrow f(x), v^k \in \hat{\partial}F(x^k) \rightarrow v \right\},$$

$\text{dom } \partial F := \{x \in \mathbb{R}^d \mid \partial F(x) \neq \emptyset\}$ and the graph of ∂F is defined as

$$\text{graph}(\partial F) := \{(x, u) \in \mathbb{R}^d \times \mathbb{R}^d \mid u \in \partial F(x)\}.$$

Fact 18. (From Rockafellar and Wets, 2009) Let $F : \mathbb{R}^d \mapsto \mathbb{R}$ be a lower semi-continuous function and consider a sequence $\{(x_k, u_k)\}_{k \geq 0} \in \text{graph}(\partial F)^{\mathbb{N}}$ such that the sequence $\{(x_k, u_k, F(x_k))\}_{k \geq 0}$ converges to a point $\{(x, u, F(x))\}$. Then $(x, u) \in \text{graph}(\partial F)$.

Definition 19. (From Attouch et al., 2013) The function $F : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is said to have the Kurdyka-Lojasiewicz property at $x^* \in \text{dom } \partial F$ if there exists $\eta \in (0, +\infty]$, a neighborhood \mathcal{U} of x^* and a continuous concave function $\varphi : [0, \eta) \rightarrow \mathbb{R}_+$ such that :

1. $\varphi(0) = 0$,
2. φ is \mathcal{C}^1 on $(0, \eta)$ and continuous at 0,
3. for all s in $(0, \eta)$, $\varphi'(s) > 0$,
4. for all $x \in \mathcal{U} \cap \{y \mid F(x^*) < F(y) < F(x^*) + \eta\}$, the Kurdyka-Lojasiewicz inequality holds

$$\varphi'(F(x) - F(x^*)) \text{dist}(0, \partial F(x)) \geq 1.$$

Proper lower semi-continuous functions which satisfy the Kurdyka-Lojasiewicz inequality at each point of $\text{dom } \partial F$ are called *KL functions*. Besides, *KL with exponent α* means the *KL property with a function $\varphi : s \mapsto cs^{1-\alpha}$* where $c > 0$. We denote this property *KL- α* .

Lemma 20. (From Bolte et al., 2007) Let $F : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ be a subanalytic function with closed domain and assume that $F|_{\text{dom } F}$ is continuous. Then for any $x \in \text{dom } F$, there exist a neighborhood $\mathcal{V} \subset \mathbb{R}^d$ of x , an exponent $\theta \in [0, 1)$ and a constant $C > 0$ such that for all $y \in \mathcal{V}$, we have

$$|F(y) - F(x)|^\theta \leq C \text{dist}(0, \partial F(y)).$$

Note that norms and in particular the Frobenius norm, the trace-norm and the group-Lasso norm satisfy the *KL property*, so the functions that we consider in this paper satisfy this property.

B.5 Critical and KW-stationary points

Definition 21. We say that $x \in \mathbb{R}^d$ is a *critical point* of F if $0 \in \partial F(x)$ where $\partial F(x)$ is defined in Definition 17.

Definition 22. Given a function $F := f_1 - f_2 + \lambda h$ where f_1 is differentiable while f_2 and h are proper, lower semi-continuous and convex, we say that $x \in \mathbb{R}^d$ is a *KW-stationary point* if there exists $u(x) \in \partial f_2(x)$ and $v(x) \in \partial h(x)$ such that

$$\nabla f_1(x) - u(x) + v(x) = 0.$$

Remark 23. Note that the Definition 21 of critical points and the Definition 22 of KW-stationary points coincide when the function f_2 is differentiable.

C The Orthogonal Procrustes Problem

Given a matrix $M \in \mathbb{R}^{p,k}$ with $p \geq k$, we use at several points in the paper the following results that were presented in the Proof of Lemma 6 in (Ge et al., 2017).

Fact 24. If $M = M_1^T M_2$, then

$$\max_{V \in \mathbb{R}^{p,k}; V^T V = I_k} \langle M, V \rangle \text{ has the same set of optima as } \min_{V \in \mathbb{R}^{p,k}; V^T V = I_k} \frac{1}{2} \|M_2 - M_1 V\|_F^2.$$

Fact 25. The optimal value of the following orthogonal Procrustes problem is given by

$$\max_{V \in \mathbb{R}^{p,k}; V^T V = I_k} \langle M, V \rangle = \|M\|_*.$$

Fact 26. If $R_1 \Sigma R_2^T$ is a complete singular value decomposition of M where $R_1 \in \mathbb{R}^{p,p}$ is such that $R_1^T R_1 = I_p$, $\Sigma \in \mathbb{R}_+^{p,k}$ has non-zero elements $\sigma_1 \geq \dots \geq \sigma_k \geq 0$ only on the diagonal and $R_2 \in \mathbb{R}^{k,k}$ is such that $R_2^T R_2 = I_k$, then an optimal solution of the orthogonal Procrustes problem is given by

$$R_1 \begin{bmatrix} I_k \\ 0_{p-k,k} \end{bmatrix} R_2^T \in \operatorname{argmax}_{V \in \mathbb{R}^{p,k}: V^T V = I_k} \langle M, V \rangle.$$

Fact 27. With the same notations as in Fact 26, if M is full-rank then, although R_1 and R_2 are not uniquely defined, the following Procrustes problem has a unique solution :

$$\operatorname{argmax}_{V \in \mathbb{R}^{p,k}: V^T V = I_k} \langle M, V \rangle = \left\{ R_1 \begin{bmatrix} I_k \\ 0_{p-k,k} \end{bmatrix} R_2^T \right\}.$$

Fact 28. If $p = k$ then $I_r \in \operatorname{argmax}_{V \in \mathbb{R}^{p,p}: V^T V = I_p} \langle M, V \rangle$ if and only if M is positive-semidefinite.

Proof. Fact 24 comes by seeing that for any $V \in \mathbb{R}^{p,k}$ such that $V^T V = I_k$, we have

$$\begin{aligned} \frac{1}{2} \|M_2 - M_1 V\|_F^2 &= \frac{1}{2} \|M_2\|_F^2 + \frac{1}{2} \|M_1 V\|_F^2 - 2 \langle M_2, M_1 V \rangle \\ &= \frac{1}{2} \|M_2\|_F^2 + \frac{1}{2} \|M_1\|_F^2 - 2 \langle M_1^T M_2, V \rangle. \end{aligned}$$

To prove Fact 25 and Fact 26, let $R_1 \Sigma R_2^T$ be a singular value decomposition of M where $R_1 \in \mathbb{R}^{p,k}$ is such that $R_1^T R_1 = I_k$, $\Sigma \in \mathbb{R}_+^{k,k}$ has nonzero elements $\sigma_1 \geq \dots \geq \sigma_k \geq 0$ only on the diagonal and $R_2 \in \mathbb{R}^{k,k}$ is such that $R_2^T R_2 = I_k$. Also, let $R_1^\perp \in \mathbb{R}^{p,p-k}$ such that $R := \begin{bmatrix} R_1 & R_1^\perp \end{bmatrix}$ satisfies $R^T R = I_p$. Writing $M = R_1 \Sigma R_2^T$ and using the change of variables $V = R_1 A R_2^T + R_1^\perp B R_2^T$, we have

$$\begin{aligned} & \max_{V \in \mathbb{R}^{p,k}: V^T V = I_k} \langle M, V \rangle & (10) \\ &= \max_{A \in \mathbb{R}^{k,k}, B \in \mathbb{R}^{p-k,k}: A^T A + B^T B = I_k} \langle R_1 \Sigma R_2^T, R_1 A R_2^T + R_1^\perp B R_2^T \rangle \\ &= \max_{A \in \mathbb{R}^{p,p}, B \in \mathbb{R}^{p-k,k}: A^T A + B^T B = I_k} \left\langle \begin{bmatrix} \Sigma \\ 0_{p-k,k} \end{bmatrix}, \begin{bmatrix} A \\ B \end{bmatrix} \right\rangle \\ &= \max_{C \in \mathbb{R}^{p,k}: C^T C = I_k} \left\langle \begin{bmatrix} \Sigma \\ 0_{p-k,k} \end{bmatrix}, C \right\rangle. \end{aligned}$$

Let $C \in \mathbb{R}^{p,k}$ such that $C^T C = I_k$, we have

$$\begin{aligned} \left\langle \begin{bmatrix} \Sigma \\ 0_{p-k,k} \end{bmatrix}, C \right\rangle &= \sum_{i=1}^k \sigma_i C_{i,i} \\ &\leq \sum_{i=1}^k \sigma_i & (11) \\ &= \|\Sigma\|_* \\ &= \|M\|_*. \end{aligned}$$

We have Inequality (11) since Σ has only nonnegative coefficients and the columns of C have unit norm so $C_{i,i} \leq 1$ for all $1 \leq i \leq k$. Besides, Inequality (11) is attained for $C = \begin{bmatrix} I_k \\ 0_{p-k,k} \end{bmatrix}$ which corresponds in

Problem (10) to $V = R_1 \begin{bmatrix} I_k \\ 0_{p-k,k} \end{bmatrix} R_2^T$. This proves Fact 25 and Fact 26.

To prove Fact 27, that is to say that $\operatorname{argmax}_{V \in \mathbb{R}^{p,k}: V^T V = I_k} \langle M, V \rangle$ is a singleton if M is full-rank, it is sufficient to notice that Inequality (11) is strict if all the σ_i are non-zero and $C_{i,i} \neq 1$ for some $1 \leq i \leq k$.

To prove Fact 28, note that $I_r \in \operatorname{argmax}_{V \in \mathbb{R}^{p,p}: V^T V = I_p} \langle M, V \rangle$ implies $\operatorname{tr}(M) = \|M\|_*$ with Fact 25 and this is only true for positive-semidefinite matrices. Conversely, if M is positive-semidefinite, then by Fact 26, $I_r \in \operatorname{argmax}_{V \in \mathbb{R}^{p,p}: V^T V = I_p} \langle M, V \rangle$. \square

D The Forward-Backward Descent Algorithm 1

Given $U \in \mathbb{R}^{p,r}$, we recall that we compute the forward direction for Algorithm 1 with the gradient $X^T XU$ of $U' \mapsto \frac{1}{2}\|XU'\|_F^2$ and z_U a subgradient of $U' \mapsto \|Y^T XU'\|_*$ whose computation is detailed in Appendix D.1.2. Setting with a slight abuse of notation $\nabla f(U) := X^T XU - z_U$, t and U_+ are then obtained with Algorithm 2 such that the (LS) condition $\tilde{F}_{t,U}^\lambda(U^+) \geq F^\lambda(U^+)$ is satisfied where

$$U_+ = \operatorname{argmin}_{U' \in \mathbb{R}^{p,r}} f(U) + \langle \nabla f(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \|U'\|_{1,2}. \quad (12)$$

D.1 Subgradients for the descent direction

If we strictly applied the subgradient-type algorithm proposed by Khamaru and Wainwright (2018) and computed a forward direction for Algorithm 1 with Fact 29, we could only prove global convergence to a KW-stationary point. Instead, we introduce in Appendix D.1.2 an additional condition on the subgradient that is leveraged in Appendix F to guarantee convergence to a critical point.

D.1.1 Subgradients of $U \mapsto \|Y^T XU\|_*$.

Thanks to Fact 25 and Fact 26, we can easily compute subgradients of $f_2 : U \mapsto \|Y^T XU\|_*$.

Fact 29. *Let $n, p \geq 0$, $r \leq \min(n, p)$, $X \in \mathbb{R}^{n,p}$, $Y \in \mathbb{R}^{n,k}$, $U \in \mathbb{R}^{p,r}$ and $R_1 DR_2^T$ be a singular value decomposition of $Y^T XU$ with $R_1 \in \mathbb{R}^{k,r}$, $R_1^T R_1 = I_r$, $D \in \mathbb{R}^{r,r}$ a diagonal matrix with nonnegative coefficients, $R_2 \in \mathbb{R}^{r,r}$ and $R_2^T R_2 = I_r$. We denote $V = R_1 R_2^T \in \mathbb{R}^{k,r}$. For any $U' \in \mathbb{R}^{p,r}$, we have*

$$\|Y^T XU'\|_* \geq \|Y^T XU\|_* + \langle X^T YV, U' - U \rangle.$$

Therefore, $X^T YV$ is a subgradient of $f_2 : U' \mapsto \|Y^T XU'\|_*$ at U .

Proof. Let $U \in \mathbb{R}^{p,r}$ and $V \in \mathbb{R}^{k,r}$ be defined as in Fact 29. Since $V^T V = R_2 R_1^T R_1 R_2^T = I_r$, we have by Fact 25 and Fact 26 :

$$\|Y^T XU\|_* = \langle V, Y^T XU \rangle. \quad (13)$$

By Fact 25, we also have for any $U' \in \mathbb{R}^{p,r}$,

$$\|Y^T XU'\|_* \geq \langle V, Y^T XU' \rangle. \quad (14)$$

Combining Equation (13) and Equation (14), we obtain

$$\|Y^T XU'\|_* \geq \|Y^T XU\|_* + \langle V, Y^T X(U' - U) \rangle.$$

□

Remark 30. *We could also obtain subgradients of f_2 using Danskin's Theorem (Danskin, 1967) but the proposed analysis in the proof of Fact 29 seems more explicit. Besides, the choice of a specific subgradient in Lemma 32 is pivotal for the global convergence analysis in Appendix F, as explained in Remark 31.*

D.1.2 Computations of z_U for Algorithm 1.

Here, we present how, given $U \in \mathbb{R}^{p,r}$, the subgradient of $f_2 : U \mapsto \|Y^T XU\|_*$ is built for Algorithm 1 and we do not assume necessarily that $X^T X$ is full-rank. Therefore, we denote $(X^T X)^{\frac{1}{2}}$ a square-root of the pseudo-inverse of $X^T X$ and, PSQ^T the reduced singular value decomposition of $(X^T X)^{\frac{1}{2}} X^T Y$. If the latter has rank ℓ then $P \in \mathbb{R}^{p,\ell}$ and $Q \in \mathbb{R}^{k,\ell}$ have orthonormal columns and $S \in \mathbb{R}^{\ell,\ell}$ is the diagonal matrix with singular values $s_1 \geq \dots \geq s_\ell > 0$. We also denote $M \in \mathbb{R}^{k,r}$ a matrix whose columns are orthonormal and belong to $\operatorname{Im} Y^T X (X^T X)^{\frac{1}{2}}$, we compute this matrix only once at the beginning of Algorithm 1 with a

Gram-Schmidt process. When $X^T X$ is invertible, the computational cost is significantly reduced since we then have $\text{Im } Y^T X (X^T X)^{\frac{1}{2}} = \text{Im } Y^T X$.

To compute z_U for Algorithm 1 - given $U \in \mathbb{R}^{p,r}$ - we first compute a singular value decomposition LDR_2^T of $Y^T XU$ with $c = \text{rank}(Y^T XU)$, $L \in \mathbb{R}^{k,r}$, $L^T L = I_r$, $D \in \mathbb{R}^{r,r}$ a diagonal matrix with nonnegative coefficients, $R_2 \in \mathbb{R}^{r,r}$ and $R_2^T R_2 = I_r$. The computational cost is $O(kr^2)$ and we write

$$\begin{aligned} L &= [L^{>0} \quad L^0], \quad \text{with } L^{>0} \in \mathbb{R}^{k,c}, L^0 \in \mathbb{R}^{k,r-c}, \\ D &= \begin{bmatrix} D^{>0} & 0_{c,r-c} \\ 0_{r-c,c} & 0_{r-c,r-c} \end{bmatrix}, \quad \text{with } D^{>0} \in \mathbb{R}^{c,c}, \\ R_2 &= [R_2^{>0} \quad R_2^0], \quad \text{with } R_2^{>0} \in \mathbb{R}^{r,c}, R_2^0 \in \mathbb{R}^{r,r-c}, \end{aligned}$$

so that

$$Y^T XU = LDR_2^T = [L^{>0} \quad L^0] \begin{bmatrix} D^{>0} & 0_{c,r-c} \\ 0_{r-c,c} & 0_{r-c,r-c} \end{bmatrix} \begin{bmatrix} R_2^{>0,T} \\ R_2^{0,T} \end{bmatrix}.$$

Clearly, the columns of $L^{>0}$ are in $\text{Im } Y^T X$ since $D^{>0} R_2^{>0} \in \mathbb{R}^{c,r}$ is full-rank. Then we apply the Gram-Schmidt process to the columns of the matrix

$$[L^{>0} \quad M] \in \mathbb{R}^{k,c+r}$$

starting from the first column of M and until we obtain $r - c$ new orthogonal vectors. The computational cost is again $O(kr^2)$. Extracting these $r - c$ vectors and denoting $\bar{L} \in \mathbb{R}^{k,r-c}$ the matrix obtained by concatenation, we define $R_1 := [L^{>0} \quad \bar{L}] \in \mathbb{R}^{k,r}$ and

$$Y^T XU = R_1 DR_2^T = [L^{>0} \quad \bar{L}] \begin{bmatrix} D^{>0} & 0_{c,r-c} \\ 0_{r-c,c} & 0_{r-c,r-c} \end{bmatrix} \begin{bmatrix} R_2^{>0,T} \\ R_2^{0,T} \end{bmatrix}.$$

Thus we obtain a singular value decomposition $R_1 DR_2^T$ of $Y^T XU$ with $\text{Im } R_1 \subset \text{Im } Y^T X (X^T X)^{\frac{1}{2}}$ at a computational cost of $O(kr^2)$. Eventually, given $U \in \mathbb{R}^{p,r}$, the subgradient of $U' \mapsto \|Y^T XU\|_*$ at U that we choose for Algorithm 1 is:

$$z_U = X^T Y R_1 R_2^T. \quad (15)$$

Remark 31. *In this paper, the condition $\text{Im } R_1 \subset \text{Im } Y^T X (X^T X)^{\frac{1}{2}}$ is only used in Lemma 32 to guarantee that $z_U \in \partial(-f_2)(U)$ where $f_2 : U' \mapsto \|Y^T XU'\|_*$. This property can then be leveraged to prove global convergence for RRR and SRRR of the iterates produced by Algorithm 1 to a critical point in the sense of Definition 21. If we do not impose this extra condition and compute a subgradient as in Fact 29, all the results still hold except for the fact that we only guarantee global convergence to a KW-stationary point in the sense of Definition 22. When $X^T X$ is invertible, we have shown that the induced computations have the same complexity $O(kr^2)$ as the computation of the SVD of $Y^T XU$.*

Lemma 32. *Given $U \in \mathbb{R}^{p,r}$ let $R_1 DR_2^T$ be a singular value decomposition of $Y^T XU$ with $R_1 \in \mathbb{R}^{k,r}$, $R_1^T R_1 = I_r$, $\text{Im } R_1 \subset \text{Im } Y^T X (X^T X)^{\frac{1}{2}}$, $D \in \mathbb{R}^{r,r}$ a diagonal matrix with nonnegative coefficients, $R_2 \in \mathbb{R}^{r,r}$ and $R_2^T R_2 = I_r$. The matrix $-z_U := -X^T Y R_1 R_2^T$ belongs to the limiting subdifferential presented in Definition 17 of the concave function $U' \mapsto -\|Y^T XU'\|_*$.*

Proof. First, with the notations of Lemma 32 and Proposition 6 of (Grave et al., 2011) that is recalled in Proposition 66, we know that when $Y^T XU$ is full-rank, the function $f_2 : U' \mapsto \|Y^T XU'\|_*$ is differentiable at U with gradient $X^T Y R_1 R_2^T$ so $-X^T Y R_1 R_2^T \in \partial(-f_2)(U)$.

Secondly, we assume that $Y^T XU$ has rank $c < r$. To prove that $-X^T Y R_1 R_2^T \in \partial(-f_2)(U)$, we exhibit a sequence $(U_k)_{k \geq 0} \in (\mathbb{R}^{p,r})^{\mathbb{N}}$ such that, as in Definition 17,

$$U^k \rightarrow U, \|Y^t XU^k\|_* \rightarrow \|Y^T XU\|_*, \text{ and } X^T Y R_1 R_2^T \in \hat{\partial}(-f_2)(U^k), \quad (16)$$

where $\hat{\partial}(-f_2)$ is the Fréchet subdifferential presented in Definition 17. Indeed, for $\epsilon > 0$, consider

$$U_\epsilon := U + \epsilon(X^T X)^{\frac{1}{2}} P S^{-1} Q^T R_1 R_2^T,$$

where PSQ^T is the reduced singular value decomposition of $(X^T X)^{\frac{1}{2}} X^T Y$. We have

$$\begin{aligned} Y^T X U_\epsilon &= Y^T X U + \epsilon Y^T X (X^T X)^{\frac{1}{2}} P S^{-1} Q^T R_1 R_2^T \\ &= R_1 D R_2^T + \epsilon Q Q^T R_1 R_2^T \\ &= R_1 D R_2^T + \epsilon R_1 R_2^T \\ &= R_1 (D + \epsilon I_r) R_2^T. \end{aligned} \tag{17}$$

Equation (17) is due to the fact that $Q Q^T R_1 = R_1$ since we assumed that $\text{Im } R_1 \subset \text{Im } Y^T X (X^T X)^{\frac{1}{2}}$ and the columns of Q form an orthonormal basis of $\text{Im } Y^T X (X^T X)^{\frac{1}{2}}$. The trace norm is therefore differentiable at $Y^T X U_\epsilon$ that is full-rank and the gradient of $U' \mapsto \|Y^T X U'\|_*$ at U_ϵ is $X^T Y R_1 R_2^T$. Defining $U_k := U_{\frac{1}{k}}$ for all $k > 0$ leads to (16). \square

D.2 The proximal operator of the group-Lasso norm

In order to highlight the fact that U_+ is simply obtained by computing ∇f and the proximal operator of the group-Lasso norm, we could equivalently write Equation (12) as

$$U_+ = \underset{U' \in \mathbb{R}^{p,r}}{\text{argmin}} \frac{1}{2} \|U' - (U - t \nabla f(U))\|_F^2 + \lambda t \|U'\|_{1,2}. \tag{18}$$

An explicit form of this proximal operator is for instance given in Equation (3.7) in (Bach et al., 2012). Given $1 \leq i \leq p$, let $[U_+]_{i,:}$ and $[U - t \nabla f(U)]_{i,:}$ denote the i -th lines of the matrices U_+ and $U - t \nabla f(U)$ respectively. Assume that $[U - t \nabla f(U)]_{i,:} \neq 0$, then we have

$$[U_+]_{i,:} = \max \left(0, 1 - \frac{\lambda t}{\|[U - t \nabla f(U)]_{i,:}\|_2} \right) [U - t \nabla f(U)]_{i,:}.$$

E The Line Search Procedure in Algorithm 2

Given $t > 0$ and $U \in \mathbb{R}^{p,r}$, we recall the definitions of $\tilde{f}_{t,U}$, $\tilde{F}_{t,U}^\lambda$ and $\gamma_t(U)$:

$$\begin{aligned} \tilde{f}_{t,U}(U') &= f(U) + \langle \nabla f(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2, \\ \tilde{F}_{t,U}^\lambda(U') &= \tilde{f}_{t,U}(U') + \lambda \|U'\|_{1,2}, \end{aligned} \tag{19}$$

$$\gamma_t(U) = -\frac{1}{t} \min_{U' \in \mathbb{R}^d} \left[\tilde{F}_{t,U}^\lambda(U') - F^\lambda(U) \right]. \tag{20}$$

E.1 A lower-bound for the decrease in terms of function values

As announced in Section 5.3, we prove that $t\gamma_t(U)$ is a lower bound for the decrease at each iteration in terms of function values.

Fact 33. *Given $U \in \mathbb{R}^{p,r}$, t and U_+ obtained with Algorithm 2, the quantity $t\gamma_t(U)$ is a lower bound for the decrease in terms of function values from U to U_+ :*

$$t\gamma_t(U) \leq F^\lambda(U) - F^\lambda(U_+).$$

Proof. Indeed, we have

$$t\gamma_t(U) = - \min_{U' \in \mathbb{R}^{p,r}} \left[\tilde{F}_{t,U}^\lambda(U') - F^\lambda(U) \right] \quad (21)$$

$$= F^\lambda(U) - \tilde{F}_{t,U}^\lambda(U_+) \quad (22)$$

$$\leq F^\lambda(U) - F^\lambda(U_+). \quad (23)$$

□

Equation (21) comes from the definition of γ_t in Equation (20). Equation (22) follows from the definition of U_+ in Equation (18). We have Equation (23) since the (LS) condition $\tilde{F}_{t,U}^\lambda(U_+) \geq F^\lambda(U_+)$ is satisfied for t and U_+ .

E.2 A lower bound on the step size with the Line Search Procedure

In this section, we prove two additional results : that the (LS) condition is satisfied as soon as $t \leq \frac{1}{L_X}$ and, that there exists $\bar{k} \in \mathbb{N}$ such that for all $k \geq \bar{k}$, we have $t_k > \frac{\beta}{L_X}$.

Lemma 34. *Let $L_X > 0$ be the largest eigenvalue of $X^T X$. For any $t \leq \frac{1}{L_X}$ and $U, U' \in \mathbb{R}^{p,r}$, we have*

$$f(U') + \lambda \|U'\|_{1,2} \leq f(U) + \langle \nabla f(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \|U'\|_{1,2} \quad (24)$$

where $y_U := X^T X U$ is the gradient of $U' \mapsto \frac{1}{2} \|XU'\|_F^2$, z_U is any subgradient of $U' \mapsto \|Y^T X U'\|_*$ and, with a slight abuse of notation, $\nabla f(U) := y_U - z_U$. Equivalently, for any $t \leq \frac{1}{L_X}$, the (LS) condition is satisfied i.e. we have

$$F^\lambda(U') \leq \tilde{F}_{t,U}^\lambda(U'). \quad (25)$$

In particular, Lemma 34 implies that Algorithm 2 terminates. This is illustrated in Figure 4.

Proof. Let $U \in \mathbb{R}^{p,r}$. On the one hand, we have for all $U' \in \mathbb{R}^{p,r}$,

$$\begin{aligned} \frac{1}{2} \|XU'\|_F^2 &= \frac{1}{2} \|X(U + (U' - U))\|_F^2 \\ &\leq \frac{1}{2} \|XU\|_F^2 + \langle X^T X U, U' - U \rangle + \frac{1}{2} \|X(U' - U)\|_F^2 \\ &\leq \frac{1}{2} \|XU\|_F^2 + \langle X^T X U, U' - U \rangle + \frac{L_X}{2} \|U' - U\|_F^2, \end{aligned} \quad (26)$$

since $L_X > 0$ is the largest eigenvalue of $X^T X$. On the other hand, since z_U is a subgradient of $U' \mapsto \|Y^T X U'\|_*$, we have for any $U' \in \mathbb{R}^{p,r}$,

$$- \|Y^T X U'\|_* \leq - \|Y^T X U\|_* - \langle z_U, U' - U \rangle. \quad (27)$$

Summing Equation(26) and Equation(27), we obtain

$$f(U') \leq \tilde{f}_{\frac{1}{L_X}, U}(U').$$

Additionally, for any $0 < t \leq \frac{1}{L_X}$, we have

$$\tilde{f}_{\frac{1}{L_X}, U}(U') \leq \tilde{f}_{t, U}(U').$$

Consequently, for any $U, U' \in \mathbb{R}^{p,r}$ and $0 < t \leq \frac{1}{L_X}$, we have

$$F^\lambda(U') = f(U') + \lambda \|U'\|_{1,2} \leq \tilde{f}_{t, U}(U') + \lambda \|U'\|_{1,2} = \tilde{F}_{t,U}^\lambda(U')$$

which is the (LS) condition. □

Consequently, if $t_k > \frac{\beta}{L_X}$ for a given $k \geq 0$, then for all $k' \geq k$ we have $t_{k'} > \frac{\beta}{L_X}$. Thus, if $\bar{t} > \frac{\beta}{L_X}$ then for all $k \geq 0$, we have $t_k > \frac{\beta}{L_X}$. If $\bar{t} \leq \frac{\beta}{L_X}$, the algorithm progressively increases the value of t and after a few first iterations, say k , we have $t_k > \frac{\beta}{L_X}$: the step size t will be larger than $\frac{\beta}{L_X}$ after a number of steps which is finite in expectation. \square

F Study of the global convergence

Khamaru and Wainwright (2018) study the convergence of subgradient-type algorithms to KW-stationary points (see Definition 22) of non-convex and non-smooth functions that can be written as a sum of three terms $F = f_1 - f_2 + \lambda h$ where f_1 is a smooth function, f_2 is a continuous and convex function, h is a possibly non-smooth, convex penalty and $\lambda \geq 0$. Some of their results can be adapted to (RRR) and (SRRR) by taking

$$\begin{aligned} f_1(U) &:= \frac{1}{2} \|XU\|_F^2, \\ f_2(U) &:= \|Y^T XU\|_*, \\ \text{and } h(U) &:= \|U\|_{1,2}. \end{aligned}$$

First, we introduce the following results by Khamaru and Wainwright (2018) that we invoke in Section F.1 and Section F.2 .

Lemma 36. (From Lemma 5 in Khamaru and Wainwright, 2018) *Let $\lambda \geq 0$ and $(U_k)_{k \geq 0}$ be the sequence generated by Algorithm 1 and $(z_k)_{k \geq 0}$ the corresponding sequence of subgradients of f_2 . For all $k \geq 0$, there is a subgradient s_{k+1} of $U \mapsto \|U\|_{1,2}$ at U_{k+1} such that*

$$U_{k+1} = U_k - t_k [\nabla f_1(U_k) - z_k + \lambda s_{k+1}], \quad (28)$$

$$F^\lambda(U_k) - F^\lambda(U_{k+1}) \geq \frac{1}{2t_k} \|U_{k+1} - U_k\|_F^2. \quad (29)$$

Furthermore, for any convergent subsequence $(U_{k_j})_{j \geq 0}$ of the sequence $(U_k)_{k \geq 0}$ with $U_{k_j} \rightarrow \bar{U}$, we have

$$\lim_{j \rightarrow +\infty} \|U_{k_j+1}\|_{1,2} = \|\bar{U}\|_{1,2}. \quad (30)$$

Lemma 36 is due to the choice of the forward-backward Algorithm 1 while the following Lemma 37 comes from the property of subanalytic functions (Bolte et al., 2007; Attouch et al., 2010, and references therein) given by Lemma 20. Indeed, norms and in particular the Frobenius norm, the trace-norm and the group-Lasso norm are subanalytic so the functions f and F^λ that we consider are subanalytic.

Lemma 37. (From Lemma 6 in Khamaru and Wainwright, 2018) *Let $\lambda \geq 0$, $(U_k)_{k \geq 0}$ be the sequence generated by Algorithm 1 and $(z_k)_{k \geq 0}$ the corresponding sequence of subgradients of f_2 . The function F^λ is constant on the set of limit points $\bar{\mathcal{U}}$ of the sequence $(U_k)_{k \geq 0}$. We denote \bar{F}^λ this limit. If we assume that $\bar{\mathcal{U}}$ contains only critical points of F^λ , then there exists constants $\theta \in [0, 1)$, $C > 0$ and $k_1 \in \mathbb{N}$ such that for all $k \geq k_1$, we have*

$$|F^\lambda(U_k) - \bar{F}^\lambda|^\theta \leq C \text{dist}(0, \nabla f_1(U_k) - z_{U_k} + \lambda \partial \|\cdot\|_{1,2}(U_k)). \quad (31)$$

F.1 Global convergence to a critical point with Algorithm 1 for RRR

The function $U \mapsto \frac{1}{2} \|XU\|_F^2$ is continuously differentiable and L_X -smooth where L_X is the largest eigenvalue of $X^T X$. The function $U \mapsto \|Y^T XU\|_*$ is continuous and convex and the difference $f(U) = \frac{1}{2} \|XU\|_F^2 - \|Y^T XU\|_*$ is bounded below by $-\frac{1}{2} \|Y\|_F^2$, indeed we have used in Section 3.1 the fact that for any $U \in \mathbb{R}^{p,r}$, we have

$$\frac{1}{2} \|XU\|_F^2 - \|Y^T XU\|_* + \frac{1}{2} \|Y\|_F^2 = \min_{V \in \mathbb{R}^{k,r}; V^T V = I_r} \frac{1}{2} \|Y - XUV^T\|_F^2 \geq 0.$$

Besides, f satisfies the Kurdyka-Łojasiewicz property, presented in Definition 19, since it is the difference of two semi-algebraic functions. Therefore, our setting satisfies the conditions of Theorem 1 and Theorem 3 in Khamaru and Wainwright (2018) and we can prove that Algorithm 1 converges to a critical point from any initial point.

F.1.1 Limit points are critical points

The following result, whose proof is inspired from Theorem 1 by Khamaru and Wainwright (2018), ensures that any limit point \bar{U} of the sequence generated by Algorithm 1 for RRR satisfies $0 \in \partial f(\bar{U})$.

Theorem 38. *Let $(U_k)_{k \geq 0}$ be the sequence generated by Algorithm 1 with $\lambda = 0$. The sequence of function values is decreasing and convergent. Besides, any limit point is a critical point of the function f .*

Proof. Equation (29) guarantees that the sequence of function values is decreasing. Since f has a finite lower-bound, the sequence of function values is convergent. Additionally, the iterates are bounded since the function is coercive *i.e.* $f(U) \rightarrow +\infty$ if $\|U\|_F \rightarrow \infty$.

To establish that the limit points are critical, consider a subsequence $(U_{k_j})_{j \geq 0}$ that converges to \bar{U} and let $(z_{k_j})_{j \geq 0}$ be the associated subsequence of subgradients. Since the sequence $(U_{k_j})_{j \geq 0}$ converges to \bar{U} , we must have by Equation (28), $\|\nabla f_1(U_{k_j}) - z_{k_j}\|_F \rightarrow 0$. The function $f_1 : U \mapsto \frac{1}{2}\|XU\|_F^2$ being continuously differentiable, we have $\nabla f_1(U_{k_j}) \rightarrow \nabla f_1(\bar{U})$ and consequently $z_{k_j} \rightarrow \bar{z} := \nabla f_1(\bar{U})$. Besides, we know by Lemma 32 that for any $j \geq 0$, we have $-z_{k_j} \in \partial(-f_2)(U_{k_j})$.

We conclude like in the proof of Theorem 1 by Khamaru and Wainwright (2018), using the graph continuity of limiting subdifferentials which we recall in Fact 18, that $-\bar{z} \in \partial(-f_2)(\bar{U})$ and $\nabla f_1(\bar{U}) - \bar{z} = 0$, meaning that $0 \in \partial(f_1 - f_2)(\bar{U}) = \partial f(\bar{U})$. \square

Remark 39. *Khamaru and Wainwright (2018) proved in an abstract but similar framework that the limit points are KW-stationary point in the sense of Definition 22, meaning that they can be stationary points for Algorithm 1. Instead, Theorem 38 guarantees, more standardly, that the limit points are critical in the sense that the limiting subdifferentials at these points contain the element 0. This is permitted by Lemma 32 which we obtained by imposing the condition $\text{Im } R_1 \subset \text{Im } Y^T X$ when computing a subgradient $X^T Y R_1 R_2^T$ of $U' \mapsto -\|Y^T X U'\|_*$, where $R_1 D R_2^T$ is a singular value decomposition of $Y^T X U$. If the condition $\text{Im } R_1 \subset \text{Im } Y^T X$ was removed, exactly the same proof as for Theorem 38 would show that the limit points are KW-stationary points.*

F.1.2 Convergence for RRR of Algorithm 1

Since f satisfies the KL property, we can prove the convergence to a critical point.

Theorem 40. *(From Theorem 3 Khamaru and Wainwright, 2018) The sequence $(U_k)_{k \geq 0}$ produced by Algorithm 1 for RRR converges to a critical point.*

The proof of Theorem 40 is identical to the proof of Theorem 3 by Khamaru and Wainwright (2018). We reproduce it here for completeness.

Proof. To prove that the sequence $(U_k)_{k \geq 0}$ has a finite length *i.e.* $\sum_{k=0}^{+\infty} \|U_k - U_{k+1}\|_F < +\infty$, we use the KL property for subanalytic functions given by Lemma 37. Let $\theta \in [0, 1)$, $C > 0$, $k_1 \in \mathbb{N}$ be defined as in Lemma 37, $k \geq k_1$ and let \bar{f} denote the limit of the sequence $\{f(U_k)\}_{k \geq 0}$. We have

$$(f(U_k) - \bar{f})^{1-\theta} - (f(U_{k+1}) - \bar{f})^{1-\theta} \geq (1-\theta)(f(U_k) - \bar{f})^{-\theta} [f(U_k) - f(U_{k+1})] \quad (32)$$

$$\geq \frac{(1-\theta)}{2t_k} (|f(U_k) - \bar{f}|)^{-\theta} \|U_k - U_{k+1}\|_F^2 \quad (33)$$

$$\geq \frac{(1-\theta)}{2Ct_k \|\nabla f_1(U_k) - z_k\|_F} \|U_k - U_{k+1}\|_F^2 \quad (34)$$

$$\geq \frac{(1-\theta)}{2C} \|U_k - U_{k+1}\|_F. \quad (35)$$

Inequality (32) follows from the concavity of $t \mapsto t^{1-\theta}$ and the inequalities $f(U_k) \geq f(U_{k+1}) \geq \bar{f}$. Inequality (33) comes from Equation (29) and the fact that $\{f(U_k)\}_{k \geq 0}$ is decreasing and converges to \bar{f} . Inequality (34) comes from Lemma 37. Finally, Inequality (35) follows from Equation (28). Summing both sides of Inequality (35) from $k = k_1$ to $k = +\infty$, we obtain

$$\begin{aligned} (f(U_{k_1}) - \bar{f})^{1-\theta} &= \sum_{k=k_1}^{+\infty} (f(U_k) - \bar{f})^{1-\theta} - (f(U_{k+1}) - \bar{f})^{1-\theta} \\ &\geq \frac{(1-\theta)}{2C} \sum_{k=k_1}^{+\infty} \|U_k - U_{k+1}\|_F \end{aligned}$$

which proves the finite length property and the convergence of the sequence $(U_k)_{k \geq 0}$. With Theorem 38, we know that this limit is a critical point. \square

F.2 Global convergence to critical points with Algorithm 1 for SRRR

In this section, we justify global convergence of Algorithm 1 for SRRR to critical points and present conditions leveraged in Lemma 44 that ensure convergence to a unique point.

F.2.1 Limit points are critical points

The function f_1 is smooth and convex, the function f_2 is continuous and convex. In addition, the function $h : U \mapsto \|U\|_{1,2}$ is clearly proper, lower semi-continuous and convex and F^λ which is bounded below satisfies the KL property. Consequently, our setting for proximal gradient descent satisfies the conditions of the first part of Theorem 2 in Khamaru and Wainwright (2018) and we can adapt this result to SRRR.

Theorem 41. *Let $(U_k)_{k \geq 0}$ be the sequence generated by Algorithm 1 with $\lambda > 0$. The sequence of function values is decreasing and convergent. Besides, any limit point is a critical point of the function F^λ .*

Proof. Equation (29) guarantees that the sequence of function values is decreasing. Since F^λ has a finite lower-bound, the sequence of function values is convergent. Additionally, the iterates are bounded since the function is coercive *i.e.* $F^\lambda(U) \rightarrow +\infty$ if $\|U\|_F \rightarrow \infty$.

To establish that the limit points are critical, consider a subsequence $(U_{k_j})_{j \geq 0}$ that converges to $\bar{U} \in \mathbb{R}^{p,r}$ and let $(z_{k_j})_{j \geq 0}$ be the associated subsequence of subgradients, like in Equation (15). Since the sequence $(U_{k_j})_{j \geq 0}$ converges to \bar{U} and f_2 is continuous, the sequence $\{f_2(U_{k_j})\}_{j \geq 0}$ converges to $f_2(\bar{U})$. Given the form of the subgradients $(z_{k_j})_{j \geq 0}$ in Equation (15), they are bounded and we can assume, passing to a subsequence if necessary, that they converge to $\bar{z} \in \mathbb{R}^d$. Besides, we know by Lemma 32 that for any $j \geq 0$, we have $-z_{k_j} \in \partial(-f_2)(U_{k_j})$. Therefore, $\{(U_{k_j}, -z_{k_j}, -f_2(U_{k_j}))\}_{j \geq 0}$ converges to $(\bar{U}, -\bar{z}, -f_2(\bar{U}))$ and, using the graph continuity of limiting subdifferentials which we recall in Fact 18, we have $-\bar{z} \in \partial(-f_2)(\bar{U})$.

We now show that $-\nabla f_1(\bar{U}) + \bar{z} \in \partial(\lambda \|\cdot\|_{1,2})(\bar{U})$. Since $(\|U_{k_j} - U_{k_j+1}\|_F)_{j \geq 0}$ converges to zero, the sequence $(U_{k_j+1})_{j \geq 0}$ converges to \bar{U} and by Equation (28), the sequence $(\|\nabla f_1(\bar{U}_{k_j}) - z_{k_j} + \lambda s_{k_j+1}\|_F)_{j \geq 0}$ also converges to zero. Since f_1 is smooth, we know that $\{\nabla f_1(U_{k_j})\}_{j \geq 0}$ converges to $\nabla f_1(\bar{U})$. Combined with the convergence of $(z_{k_j})_{j \geq 0}$ to \bar{z} , it shows that $(\lambda s_{k_j+1})_{j \geq 0}$ converges to $\bar{s} := -\nabla f_1(\bar{U}) + \bar{z}$. With Equation (30) in Lemma 36, we also have that $(\lambda \|U_{k_j+1}\|_{1,2})_{j \geq 0}$ converges to $\lambda \|\bar{U}\|_{1,2}$. All this leads to the convergence of $\{(U_{k_j+1}, \lambda s_{k_j+1}, \lambda \|U_{k_j+1}\|_{1,2})\}_{j \geq 0}$ to $(\bar{U}, \lambda \bar{s}, \lambda \|\bar{U}\|_{1,2})$. Consequently, the graph continuity in Fact 18 guarantees that $\lambda \bar{s} \in \partial(\lambda \|\cdot\|_{1,2})(\bar{U})$. Finally, we conclude that $\nabla f_1(\bar{U}) - \bar{z} + \lambda \bar{s} = 0 \in \partial F^\lambda(\bar{U})$ *i.e.* \bar{U} is a critical point of F^λ . \square

Remark 42. *The same comments as in Remark 39 hold for Theorem 41 and the comparison between its proof and the proof of Theorem 2 by Khamaru and Wainwright (2018).*

F.2.2 Convergence for SRRR of Algorithm 1

In order to prove convergence of the sequence $(U_k)_{k \geq 0}$, Theorem 4 of Khamaru and Wainwright (2018) formally requires that f_2 is a smooth function, a requirement which is not met by $U \mapsto \|Y^T X U\|_*$. Nonetheless, an inspection of the proof shows that local smoothness in a neighborhood of the critical points of the function is sufficient. More precisely, the same proof as for Theorem 4 in Khamaru and Wainwright (2018) can be used for SRRR as long as we can guarantee that there exists $k_1 \geq 0$ such that for all $k \geq k_1$, the iterates $(U_k)_{k \geq k_1}$ lie in a compact subset where f is locally smooth. We denote \bar{U} the set of limit points of the sequence $(U_k)_{k \geq 0}$ and for any $S \subset \{1, \dots, p\}$, X_S is the matrix formed by keeping the columns of X indexed by S .

Theorem 43. *Assume that*

- $\mathcal{H}1$: *The step sizes $(t_k)_{k \geq 0}$ produced by Algorithm 1 are upper bounded by a constant $d > 0$.*
- $\mathcal{H}2$: *The set of limit points \bar{U} of the sequence produced by Algorithm 1 is a subset of the local minima of F^λ and contains only matrices with at least r non-zero rows.*
- $\mathcal{H}3$: *For any $S \subset \{1, \dots, p\}$ of cardinality at least r , the matrix $X_S^T Y$ is full-rank.*

Then the sequence $(U_k)_{k \geq 0}$ produced by Algorithm 1 for SRRR converges to a critical point.

The assumptions $\mathcal{H}1$ and $\mathcal{H}2$ are used in the proof of Theorem 43. The assumption $\mathcal{H}2$ will hold unless local minima are so sparse that the number of selected variables is strictly smaller than r in which case the rank constraint becomes essentially useless. The assumption $\mathcal{H}3$ will hold with probability one if X and Y contain for example additive noise. It is leveraged in Appendix K.1 to prove Lemma 44 that we introduce below with Lemma 45 and Lemma 46 before giving the proof of Theorem 43.

Lemma 44. *With Assumption $\mathcal{H}3$, any local minimum U^* of (SRRR) which has at least r non-zero rows must be full-rank.*

Put differently, Assumption $\mathcal{H}2$ and Assumption $\mathcal{H}3$ combined with Lemma 44 imply that the set of limit points \bar{U} contains only full-rank matrices. The next lemma ensures that the function $f_2 : U \mapsto \|Y^T X U\|_*$ is differentiable at such points, it is proved in Appendix K.2.

Lemma 45. *With Assumption $\mathcal{H}3$, let U^* be a full-rank local minimum of (SRRR). Necessarily, $Y^T X U^*$ is full-rank.*

Lemma 45 is essential to prove locally a Lipschitz gradients property which is formalized in Lemma 46, proved in Appendix K.3.

Lemma 46. *With Assumption $\mathcal{H}2$ and Assumption $\mathcal{H}3$, there exists $M > 0$ and $k_1 \geq 0$ such that for any $k \geq k_1$, f is differentiable at U_k, U_{k+1} and we have*

$$\|\nabla f(U_k) - \nabla f(U_{k+1})\|_F \leq M \|U_k - U_{k+1}\|_F. \quad (36)$$

Proof of Theorem 43 . Let $k_1 \geq 0$ be defined as in Lemma 46. For $k \geq k_1$ we denote z_k a gradient of f_2 obtained through the update in Algorithm 1 and s_k a subgradient of $U \mapsto \lambda \|U\|_{1,2}$ at U_k . Let $k > k_1$, we have

$$\|\nabla f_1(U_k) - z_k + \lambda s_k\|_F = \|(\nabla f_1(U_k) - z_k) + (z_{k-1} - \nabla f_1(U_{k-1})) + \frac{1}{t_{k-1}}(U_{k-1} - U_k)\|_F \quad (37)$$

$$\leq \|(\nabla f_1(U_k) - z_k) - (\nabla f_1(U_{k-1}) - z_{k-1})\|_F + \frac{1}{t_{k-1}} \|U_{k-1} - U_k\|_F \quad (38)$$

$$\begin{aligned} &= \|\nabla f(U_k) - \nabla f(U_{k-1})\|_F + \frac{1}{t_{k-1}} \|U_{k-1} - U_k\|_F \\ &\leq (M + \frac{1}{t_{k-1}}) \|U_k - U_{k-1}\|_F. \end{aligned} \quad (39)$$

Inequality (37) follows from the update in Algorithm 1. Inequality (38) comes from the triangle inequality. Inequality (39) is due to Equation (36).

The second argument we give is similar to Equation (35) in the proof of Theorem 40. Since the functions we consider are subanalytic, we can consider $\theta \in [0, 1)$, $C > 0$ and $k_2 \geq k_1$ defined as in Lemma 37. Let \bar{F}^λ denote the limit of the sequence $\{F^\lambda(U_k)\}_{k \geq 0}$ and $k \geq k_2$, we have

$$(F^\lambda(U_k) - \bar{F}^\lambda)^{1-\theta} - (F^\lambda(U_{k+1}) - \bar{F}^\lambda)^{1-\theta} \geq (1-\theta) [F^\lambda(U_k) - \bar{F}^\lambda]^{-\theta} [F^\lambda(U_k) - F^\lambda(U_{k+1})] \quad (40)$$

$$\geq \frac{(1-\theta)}{2t_k} [[F^\lambda(U_k) - \bar{F}^\lambda]]^{-\theta} \|U_k - U_{k+1}\|_F^2 \quad (41)$$

$$\geq \frac{(1-\theta)}{2Ct_k \|\nabla f(U_k) + \lambda s_k\|_F} \|U_k - U_{k+1}\|_F^2 \quad (42)$$

$$\geq \frac{(1-\theta)}{2Cd \|\nabla f(U_k) + \lambda s_k\|_F} \|U_k - U_{k+1}\|_F^2. \quad (43)$$

Inequality (40) follows from the concavity of $t \mapsto t^{1-\theta}$ and the inequalities $F^\lambda(U_k) \geq F^\lambda(U_{k+1}) \geq \bar{F}^\lambda$. Inequality (41) comes from Equation (29) and the fact that $\{F^\lambda(U_k)\}_{k \geq 0}$ is decreasing and converges to \bar{F}^λ . Inequality (42) comes from Lemma 37. Finally, Inequality (43) follows from Assumption $\mathcal{H}1$ in Theorem 43. Combining Inequality (39) with Inequality (43), we obtain

$$\begin{aligned} (F^\lambda(U_k) - \bar{F}^\lambda)^{1-\theta} - (F^\lambda(U_{k+1}) - \bar{F}^\lambda)^{1-\theta} &\geq \frac{(1-\theta)}{2Cd(M + \frac{1}{t_{k-1}})} \frac{\|U_k - U_{k+1}\|_F^2}{\|U_{k-1} - U_k\|_F} \\ &\geq \frac{(1-\theta)}{2Cd(M + \frac{1}{\min(\frac{1}{L_X}, t_{-1})})} \frac{\|U_k - U_{k+1}\|_F^2}{\|U_{k-1} - U_k\|_F}. \end{aligned} \quad (44)$$

Equation (44) follows from Fact 35. The rest of the proof leads to the finite length property and completely follows the proof of Theorem 4 in Khamaru and Wainwright (2018) since they also leverage only the local property of Lipschitz gradients in a compact set. We denote

$$\Delta_k := C' [(F^\lambda(U_k) - \bar{F}^\lambda)^{1-\theta} - (F^\lambda(U_{k+1}) - \bar{F}^\lambda)^{1-\theta}] \quad (45)$$

$$\text{where } C' := \frac{2Cd \left[M + \max\left(\frac{L_X}{\beta}, \frac{1}{t_{-1}}\right) \right]}{(1-\theta)}.$$

Equation (44) can be rewritten

$$\|U_k - U_{k+1}\|_F \leq \sqrt{\Delta_k \|U_{k-1} - U_k\|_F}.$$

Summing from $j = k_2 + 1$ to $j = k$, we obtain

$$\begin{aligned} \sum_{j=k_2+1}^k \|U_j - U_{j+1}\|_F &\leq \sum_{j=k_2+1}^k \sqrt{\Delta_j \|U_{j-1} - U_j\|_F} \\ &\leq \sum_{j=k_2+1}^k \frac{1}{2} \Delta_j + \frac{1}{2} \|U_{j-1} - U_j\|_F \end{aligned} \quad (46)$$

$$\leq \frac{C'}{2} (F^\lambda(U_{k_2+1}) - \bar{F}^\lambda)^{1-\theta} + \frac{1}{2} \sum_{j=k_2+1}^k \|U_{j-1} - U_j\|_F. \quad (47)$$

Inequality (46) follows from the inequality of arithmetic and geometric means. Inequality (47) comes from Equation (45). Rewriting Inequality (47), we have for all $k \geq k_2 + 2$,

$$\begin{aligned} & \left[\frac{1}{2} \sum_{j=k_2+1}^{k-1} \|U_j - U_{j+1}\|_F \right] + \left[\frac{1}{2} \sum_{j=k_2+2}^k \|U_{j-1} - U_j\|_F \right] + \|U_k - U_{k+1}\|_F \\ & \leq \frac{C'}{2} (F^\lambda(U_{k_2+1}) - \bar{F}^\lambda)^{1-\theta} + \left[\frac{1}{2} \sum_{j=k_2+2}^k \|U_{j-1} - U_j\|_F \right] + \frac{1}{2} \|U_{k_2} - U_{k_2+1}\|_F. \end{aligned}$$

This last inequality implies that

$$\begin{aligned} \frac{1}{2} \sum_{j=k_2+1}^{k-1} \|U_j - U_{j+1}\|_F & \leq \frac{C'}{2} (F^\lambda(U_{k_2+1}) - \bar{F}^\lambda)^{1-\theta} + \frac{1}{2} \|U_{k_2} - U_{k_2+1}\|_F - \|U_k - U_{k+1}\|_F \\ & \leq \frac{C'}{2} (F^\lambda(U_{k_2+1}) - \bar{F}^\lambda)^{1-\theta} + \frac{1}{2} \|U_{k_2} - U_{k_2+1}\|_F \\ & < +\infty. \end{aligned}$$

Eventually, we conclude that the sequence $(U_k)_{k \geq 0}$ has finite length and therefore converges to an element $\bar{U} \in \mathbb{R}^{p,r}$. With Theorem 41, we know that \bar{U} is a critical point. \square

G Proofs for section 5.1

In this section, we are going to prove Equation (7), Lemma 1 and Lemma 2. We maintain the following assumptions :

$$r \leq \ell, \tag{48}$$

$$s_1 > \dots > s_\ell. \tag{49}$$

At first, to widen the scope of our results, we will not make the assumption

$$X^T X \text{ is invertible.} \tag{50}$$

Assumption (50) will play a key role in the analysis and impact the results. We will precise what it implies for the analysis when it is satisfied and when it is not.

G.1 Proof of Equation (7)

While we assumed that X is full-rank in the core of the article, we do not make this assumption in this section to prove a more general result than Equation (7). Of course, the latter can be obtained as a special case. Let $m \leq p$ be the rank of X and consider

$$K D^2 K^T \text{ the reduced singular value decomposition of } X^T X,$$

with $K \in \mathbb{R}^{p,m}$, $K^T K = I_m$ and $D \in \mathbb{R}^{m,m}$ a diagonal matrix with positive entries on the diagonal. We also write

$$(X^T X)^\dagger := K D^{-2} K^T \text{ the pseudo-inverse of } X^T X,$$

$$(X^T X)^{\frac{1}{2}} := K D^{-1} K^T \text{ a square-root of } (X^T X)^\dagger,$$

and

$$(X^T X)^{\frac{1}{2}} := K D K^T \text{ a square root of } X^T X.$$

Let

$$K^\perp \in \mathbb{R}^{p,p-m} \text{ such that } [K \ K^\perp]^T [K \ K^\perp] = I_p.$$

Here, we denote PSQ^T the reduced singular values of $(X^T X)^{\frac{1}{2}} X^T Y$, with $\ell := \text{rank}(X^T X)^{\frac{1}{2}} X^T Y \leq \min(m, k)$, $P \in \mathbb{R}^{p,\ell}$, $S \in \mathbb{R}^{\ell,\ell}$ and $Q \in \mathbb{R}^{\ell,k}$. We also define $P^\perp \in \mathbb{R}^{p,m-\ell}$ such that the columns of the matrix $[P \ P^\perp]$ form an orthonormal basis of $\text{Im } X^T$. If X is full-rank, this definition corresponds indeed with the matrices that were introduced in Section 5.1. The definition of τ is :

$$\tau : \begin{cases} \mathbb{R}^{\ell,r} \times \mathbb{R}^{m-\ell,r} \times \mathbb{R}^{p-m,r} \rightarrow \mathbb{R}^{p,r} \\ (A, C, N) \mapsto (X^T X)^{\frac{1}{2}} [P \ P^\perp] \begin{bmatrix} S & 0 \\ 0 & I_{m-\ell} \end{bmatrix} \begin{bmatrix} A \\ C \end{bmatrix} + K^\perp N \end{cases} \quad (51)$$

Of course, under the additional assumption that $X^T X$ is invertible, the term $K^\perp N$ would be removed and τ would be the same as the one we defined in Equation (6).

We define $f_{a,c,n} := f \circ \tau$ and we prove

$$f_{a,c,n}(A, C, N) = \frac{1}{2} \|SA\|_F^2 - \|S^2 A\|_* + \frac{1}{2} \|C\|_F^2. \quad (52)$$

Equation (7) can be obtained similarly if $X^T X$ is invertible.

Proof of Equation (52). Let $(A, C, N) \in \mathbb{R}^{\ell,r} \times \mathbb{R}^{m-\ell,r} \times \mathbb{R}^{p-m,r}$, we have

$$f_{a,c,n}(A, C, N) = f \circ \tau(A, C, N) \quad (53)$$

$$= f((X^T X)^{\frac{1}{2}}(PSA + P^\perp C) + K^\perp N) \quad (54)$$

$$= \frac{1}{2} \|(X^T X)^{\frac{1}{2}}((X^T X)^{\frac{1}{2}}(PSA + P^\perp C) + K^\perp N)\|_F^2 - \|Y^T X((X^T X)^{\frac{1}{2}}(PSA + P^\perp C) + K^\perp N)\|_* \quad (55)$$

$$= \frac{1}{2} \|(X^T X)^{\frac{1}{2}}(X^T X)^{\frac{1}{2}}(PSA + P^\perp C)\|_F^2 - \|Y^T X(X^T X)^{\frac{1}{2}}(PSA + P^\perp C)\|_* \quad (56)$$

$$= \frac{1}{2} \|PSA\|_F^2 + \frac{1}{2} \|P^\perp C\|_F^2 - \|QSP^T(PSA + P^\perp C)\|_*, \quad (57)$$

$$= \frac{1}{2} \|SA\|_F^2 + \frac{1}{2} \|C\|_F^2 - \|QS^2 A\|_* \quad (58)$$

$$= \frac{1}{2} \|SA\|_F^2 - \|S^2 A\|_* + \frac{1}{2} \|C\|_F^2. \quad (59)$$

□

Equation (53) follows from the definition of $f_{a,c,n}$ and Equation (54) from the definition of τ . Equation (55) follows from the definition of f and since for all $M \in \mathbb{R}^{p,r}$, we have $\|XM\|_F^2 = \|(X^T X)^{\frac{1}{2}} M\|_F^2$. We have Equation (56) since $XK^\perp = 0$. Equation (57) comes from the facts that $P, P^\perp \in \text{Im } X$ and $(X^T X)^{\frac{1}{2}}(X^T X)^{\frac{1}{2}}$ acts like the identity on $\text{Im } X^T$ for the first term and $QSP^T = Y^T X(X^T X)^{\frac{1}{2}}$ for the second term. We have Equation (58) because $[P \ P^\perp]$ is orthogonal and Equation (59) because the columns of Q are orthogonal.

G.2 Proof of Lemma 1

We denote Ω_a^* the set of minima of $f_a : A \in \mathbb{R}^{\ell,r} \mapsto \frac{1}{2} \|SA\|_F^2 - \|S^2 A\|_*$ where $S = \text{diag}(s_1 > \dots > s_\ell) \in \mathbb{R}^{\ell,\ell}$. To prove that $\Omega_a^* = \{\tilde{I}R \mid R \in \mathcal{O}_r\}$ with $\tilde{I} = (1_{i=j})_{1 \leq i \leq \ell, 1 \leq j \leq r} \in \mathbb{R}^{\ell,r}$, first note that the two following problems have the same optimal value :

$$\min_{A \in \mathbb{R}^{\ell,r}, V \in \mathbb{R}^{\ell,r}} f_{a,v}(A, V) \text{ where } f_{a,v}(A, V) := \frac{1}{2} \|S - SAV^T\|_F^2, \quad (60)$$

$$\min_{A \in \mathbb{R}^{\ell,r}, V \in \mathbb{R}^{\ell,r}: V^T V = I_r} f_{a,v}(A, V). \quad (61)$$

Indeed, for any $A, V \in \mathbb{R}^{\ell, r}$, there exists $A', V' \in \mathbb{R}^{\ell, r}$ such that $V^T V = I_r$ and $AV^T = A'V'^T$. For instance, the matrices can be obtained from the singular value decomposition $R_1 \Sigma R_2^T$ of AV^T by taking $A' = R_1 \Sigma$ and $V' = R_2$. Besides, given $A \in \mathbb{R}^{\ell, r}$ and $V \in \mathbb{R}^{\ell, r}$, we have

$$f_{a,v}(A, V) = \frac{1}{2} \|S - SAV^T\|_F^2 = \frac{1}{2} \|S\|_F^2 + \frac{1}{2} \|SAV^T\|_F^2 - \langle S, SAV^T \rangle.$$

Defining $V_A \in \operatorname{argmax}_{V \in \mathbb{R}^{\ell, r}; V^T V = I_r} \langle S, SAV^T \rangle$ and using Fact 25, we obtain

$$\frac{1}{2} \|S - SAV_A^T\|_F^2 = \frac{1}{2} \|S\|_F^2 + \frac{1}{2} \|SA\|_F^2 - \|S^2 A\|_*.$$

Consequently, if A is a minimizer of

$$\min_{A \in \mathbb{R}^{\ell, r}} f_a(A), \quad (62)$$

where $f_a(A) = \frac{1}{2} \|SA\|_F^2 - \|S^2 A\|_*$, then (A, V_A) is a minimizer of Problem (60). This means in particular that SAV_A^T is a minimizer of

$$\min_{M \in \mathbb{R}^{\ell, \ell}; \operatorname{rank}(M) \leq r} \frac{1}{2} \|S - M\|_F^2.$$

The matrix SAV_A^T must be equal to the best low-rank approximation for the Frobenius norm of S and, by the Eckart-Young-Mirsky theorem, this best approximation is $S\tilde{I}\tilde{I}^T$ with $\tilde{I} = \begin{bmatrix} I_r \\ 0 \end{bmatrix} \in \mathbb{R}^{\ell, r}$ since we have assumed that the values on the diagonal of S are strictly decreasing.

The matrix S is invertible so we must have $AV^T = \tilde{I}\tilde{I}^T$ which is equivalent, if $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$ and $V = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix}$ with $A_1, V_1 \in \mathbb{R}^{r, r}$ and $A_2, V_2 \in \mathbb{R}^{\ell-r, \ell-r}$, to :

$$\begin{bmatrix} A_1 V_1^T & A_1 V_2^T \\ A_2 V_1^T & A_2 V_2^T \end{bmatrix} = \begin{bmatrix} I_r & 0 \\ 0 & 0 \end{bmatrix} \quad (63)$$

The second line $A_2 V_2^T = 0$ implies $A_2 = 0$ since $V^T V = I_r$. From the first line of the matrices in Equation (63), $A_1 V_1^T = I_r$ implies that A_1 is invertible so $A_1 V_2^T = 0$ implies that $V_2 = 0$ and A_1 has to be orthogonal as it is the inverse of V_1^T . Put differently, $A^T = \begin{bmatrix} V_1^T & 0_{r, \ell-r} \end{bmatrix} = V_1^T \tilde{I}^T$ where V_1 is an orthogonal square matrix *i.e.* an orthogonal matrix. Thus, any optimum A belongs to $\Omega_a^* := \{\tilde{I}R \mid R \in \mathcal{O}_r\}$. Conversely, for any $R \in \mathcal{O}_r$ we have $f_a(\tilde{I}R) = \frac{1}{2} \|S\tilde{I}\|_F^2 - \|S^2 \tilde{I}\|_*$: this implies that all the elements in Ω_a^* are optima.

G.3 Proof of Lemma 2

We show that all local minima of f_a are global. The result is the same for f given that $f \circ \tau(A, C) = f_a(A) + \frac{1}{2} \|C\|_F^2$ and τ is the invertible linear transformation defined in Equation (51). First we start by eliminating the possibility of having a local maximum other than 0 with the following result.

Lemma 47. *Only 0 can be a local maximum of f_a .*

Proof. For any A , the restriction of f_a to the one-dimensional set $\mathcal{D}_A := \{\alpha A, \alpha \geq 0\}$ is a convex polynomial function of degree 2. Indeed, for any $\alpha > 0$, we have

$$f_a(\alpha A) = \frac{\alpha^2}{2} \|SA\|_F^2 - \alpha \|S^2 A\|_*.$$

Since $S \in \mathbb{R}^{\ell, \ell}$ is an invertible diagonal matrix, only 0 can be a local maximum of f_a . □

Corollary 48. *The zero matrix is indeed a local maximum of the function f_a .*

Proof. Thanks to the equivalence of norms in finite dimensions and the fact that S has only positive elements on its diagonal, we know that there exists $c, d > 0$ such that for any $A \in \mathbb{R}^{\ell, r}$, $t > 0$, we have

$$f_a(0 + tA) \leq c\|A\|_F^2 t^2 - d\|A\|_F t.$$

The zero matrix is necessary a local maximum. \square

To deal with critical points, we treat separately rank-deficient matrices and full-rank matrices. The following result, proved in Appendix K.4, considers the case of rank-deficient matrices.

Lemma 49. *Let $A \in \mathbb{R}^{\ell, r}$ be a rank-deficient matrix, there exists $B \in \mathbb{R}^{\ell, r}$ such that $\|B\|_F = 1$ and $\delta > 0$ such that for all $-\delta < t < \delta$, we have*

$$f_a(A + tB) \leq f_a(A) - \frac{s_\ell^2}{2}|t|.$$

Therefore, no rank-deficient matrix can be a local minimum of f_a .

In order to deal with full-rank matrices and having already described the set of optima, we characterize the set of full-rank critical points. Consider the set \mathcal{P} of permutations $\pi : \llbracket 1; \ell \rrbracket \rightarrow \llbracket 1; \ell \rrbracket$ such that $\pi(1) < \dots < \pi(r)$ and simultaneously $\pi(r+1) < \dots < \pi(\ell)$. For any $\pi \in \mathcal{P}$, we denote

$$\Pi_\pi := (1_{i=\pi(j)})_{1 \leq i \leq \ell, 1 \leq j \leq r} \in \mathbb{R}^{\ell, r}. \quad (64)$$

Note that the sole purpose of the condition $\pi(r+1) < \dots < \pi(\ell)$ is to have a one-to-one correspondence between the set of permutations \mathcal{P} and the set of matrices $\{\Pi_\pi \mid \pi \in \mathcal{P}\}$. We have the following result, proved in Appendix K.5.

Lemma 50. *If the values of the diagonal matrix S are strictly decreasing i.e. $S = \text{diag}(s_1 > \dots > s_\ell)$, then the set Ω_a^s of differentiable critical points for problem (62) is the image by linear transformations from $\mathbb{R}^{r, r}$ to $\mathbb{R}^{\ell, r}$ of \mathcal{O}_r :*

$$\Omega_a^s = \{\Pi_\pi R \mid \pi \in \mathcal{P}, R \in \mathcal{O}_r\}$$

Besides, Ω_a^s contains only global minima and saddle points.

We could have an even more precise description of the behavior of f_a around the saddle points with Theorem 55 and Corollary 57 (given below). Saddle points are in fact strict saddle points i.e. the Hessian at these points has at least one negative eigenvalue. However, that is not necessary here.

We can now prove Lemma 2.

Proof of Lemma 2. We know from Lemma 49 that a rank-deficient matrix can not be a local minimum. The function f_a is differentiable at $A \in \mathbb{R}^{\ell, r}$ if and only if A is full-rank¹. Finally, Lemma 50 details all critical points where A is full-rank, they are either global minima or saddle points. \square

H Proofs for Section 5.2

H.1 Proof of Lemma 3

In Section 5.2, we have introduced for any $A \in \mathbb{R}^{p, r}$,

$$\Pi_{\Omega_a^*}(A) := \underset{B \in \Omega_a^*}{\operatorname{argmin}} \|B - A\|_F^2$$

and

$$\mathcal{C}_a(R) := \{A \in \mathbb{R}^{\ell, r} \mid \tilde{I}R \in \Pi_{\Omega_a^*}(A)\}. \quad (65)$$

¹Details about the derivative of the trace-norm are given in Proposition 66.

First, we prove Equation (8) that describes $\mathcal{C}_a(I_r)$. According to Lemma 1, $\Omega_a^* = \{\tilde{I}R \mid R \in \mathcal{O}_r\}$, so with Fact 24, we could have equivalently defined $\Pi_{\Omega_a^*}(A)$ with $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix}$, $A_1 \in \mathbb{R}^{r,r}$ and $A_2 \in \mathbb{R}^{\ell-r,r}$ as

$$\operatorname{argmin}_{\tilde{I}R: R \in \mathcal{O}_r} \|\tilde{I}R - A\|_F^2 = \operatorname{argmax}_{\tilde{I}R: R \in \mathcal{O}_r} \langle \tilde{I}R, A \rangle = \tilde{I} \operatorname{argmax}_{R \in \mathcal{O}_r} \langle R, \tilde{I}^T A \rangle = \tilde{I} \operatorname{argmax}_{R \in \mathcal{O}_r} \langle R, A_1 \rangle. \quad (66)$$

By Fact 28, we have that $I_r \in \operatorname{argmax}_{R: R \in \mathcal{O}_r} \langle R, A_1 \rangle$ if and only if A_1 is positive-semidefinite. This proves Equation (8).

Secondly, the equality $\mathcal{C}_a(R) = \{AR \mid A \in \mathcal{C}_a(I_r)\}$ basically stems from the definition of $\Pi_{\Omega_a^*}$ since

$$\|\tilde{I} - A\|_F^2 = \|\tilde{I}R - AR\|_F^2. \quad (67)$$

Indeed, Equation (67) implies that $A \in \mathcal{C}_a(I_r)$ if and only if $AR \in \mathcal{C}_a(R)$.

Finally, to prove that $\cup_{R \in \mathcal{O}_r} \mathcal{C}_a(R) = \mathbb{R}^{\ell,r}$, consider $M \in \mathbb{R}^{\ell,r}$ and $B_M \in \operatorname{argmin}_{B \in \Omega_a^*} \|B - M\|_F^2$. According to Lemma 1, $\Omega_a^* := \{\tilde{I}R \mid R \in \mathcal{O}_r\}$ is compact. Therefore, there exists $R \in \mathcal{O}_r$ such that $B_M = \tilde{I}R$. Obviously, the definition of $\mathcal{C}_a(R)$ given in Equation (65) implies that $M \in \mathcal{C}_a(R)$.

The following fact gives more details on the structure of the cones that we built.

Fact 51. *The relative interiors² of all the cones partition the set of matrices $[A_1^T \ A_2^T]^T$ such that $A_1 \in \mathbb{R}^{r,r}$ is invertible and $A_2 \in \mathbb{R}^{\ell-r,r}$.*

Proof. First, since the relative interior of the set \mathcal{S}_r^+ of positive-semidefinite matrices is the set \mathcal{S}_r^{++} of positive-definite matrices, given $R \in \mathcal{O}_r$, the relative interior of the cone $\mathcal{C}_a(R)$ is the set

$$\left\{ \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} R \mid A_1 \in \mathcal{S}_r^{++}, A_2 \in \mathbb{R}^{\ell-r,r} \right\}.$$

Secondly, according to Equation (66), the matrix $A = [A_1^T \ A_2^T]^T \in \mathcal{C}_a(R)$ with $R \in \mathcal{O}_r$ if and only if $R \in \operatorname{argmax}_{R' \in \mathcal{O}_r} \langle R', A_1 \rangle$. According to Fact 27, there is a unique element in $\operatorname{argmax}_{R' \in \mathcal{O}_r} \langle R', A_1 \rangle$ if A_1 has full rank. Therefore, given $[A_1^T \ A_2^T]^T$ such that $A_1 \in \mathbb{R}^{r,r}$ is invertible and $A_2 \in \mathbb{R}^{\ell-r,r}$, there is a unique $R \in \mathcal{O}_r$ such that $A \in \mathcal{C}_a(R)$. \square

H.2 Proof of Theorem 4

First, in order to simplify the computations, we introduce the change of variables $M = SA$ and the function

$$f_m : M \in \mathbb{R}^{\ell,r} \mapsto \frac{1}{2} \|M\|_F^2 - \|SM\|_*.$$

Note that for any $M \in \mathbb{R}^{\ell,r}$, we have $f_m(M) = f_a(S^{-1}M)$ and $\min_M f_m(M)$ is the form taken by (RRR) if X is the identity and $Y = S$ is a diagonal matrix.

As in section G.3, we consider the set \mathcal{P} of permutations $\pi : \llbracket 1; \ell \rrbracket \rightarrow \llbracket 1; \ell \rrbracket$ such that $\pi(1) < \dots < \pi(r)$ and simultaneously $\pi(r+1) < \dots < \pi(\ell)$. For any $\pi \in \mathcal{P}$, we denote

$$\Pi_\pi := (1_{i=\pi(j)})_{1 \leq i \leq \ell, 1 \leq j \leq r} \in \mathbb{R}^{\ell,r}. \quad (68)$$

With the proposed change of variables, the differentiable critical points of f_m are simply obtained from the critical points of f_a given in Lemma 50.

Lemma 52. *If the values of the diagonal matrix S are strictly decreasing, then the set Ω_m^s of differentiable critical points of f_m is the image by linear transformations from $\mathbb{R}^{r,r}$ to $\mathbb{R}^{\ell,r}$ of \mathcal{O}_r :*

$$\Omega_m^s = \{S\Pi_\pi R \mid \pi \in \mathcal{P}, R \in \mathcal{O}_r\}.$$

²Given a set in a Euclidean space, its relative interior is the interior of this set within the subspace spanned by its elements.

The following result describes the eigenvectors of the Hessian of f_m at a critical point $S\Pi_\pi R$. It is proved in Appendix K.6. We write $S^2 = \text{diag}(\sigma_1 > \dots > \sigma_\ell)$ with $\sigma_\ell > 0$. For $1 \leq i_0 \leq \ell$, $1 \leq j_0 \leq r$, we denote $E_{i_0, j_0} = e_{i_0} e_{j_0}^T \in \mathbb{R}^{\ell, r}$.

Theorem 53. *Let $S\Pi_\pi R$ be a differentiable critical point of f_m , with $\pi \in \mathcal{P}$ and $R \in \mathcal{O}_r$. Then f_m is twice differentiable at $S\Pi_\pi R$, let \mathcal{H}_m denote its Hessian at $S\Pi_\pi R$.*

- For $1 \leq i < j \leq r$, $S^{-1}(E_{\pi(i),j} + E_{\pi(j),i})R$ is an eigenvector of \mathcal{H}_m associated to the eigenvalue 1.
- For $1 \leq i \leq r$, $S^{-1}E_{\pi(i),i}R$ is an eigenvector of \mathcal{H}_m associated to the eigenvalue 1.
- For $1 \leq i < j \leq r$, $S(E_{\pi(i),j} - E_{\pi(j),i})R$ is an eigenvector of \mathcal{H}_m associated to the eigenvalue 0.
- For $r+1 \leq k \leq \ell$, $1 \leq j \leq r$, $E_{\pi(k),j}R$ is an eigenvector of \mathcal{H}_m associated to the eigenvalue $1 - \frac{\sigma_{\pi(k)}}{\sigma_{\pi(j)}}$.

Remark 54. *At an optimum $\tilde{S}R$ of f_m with $R \in \mathcal{O}_r$, the largest eigenvalue of the Hessian is 1 and the smallest positive eigenvalue is $1 - \frac{\sigma_{\pi(r+1)}}{\sigma_{\pi(r)}}$.*

Since we used the change of variables $M = SA$, Theorem 53 can be adapted to the function f_a .

Theorem 55. *Let $\Pi_\pi R$ be a differentiable critical point of f_a , with $\pi \in \mathcal{P}$ and $R \in \mathcal{O}_r$. Then f_a is twice differentiable at $\Pi_\pi R$, let \mathcal{H}_a denote its Hessian at $\Pi_\pi R$.*

- For $1 \leq i < j \leq r$, $(E_{\pi(i),j} + E_{\pi(j),i})R$ is an eigenvector of \mathcal{H}_a associated to the eigenvalue $(\sigma_{\pi(i)}^{-1} + \sigma_{\pi(j)}^1)^{-1}$.
- For $1 \leq i \leq r$, $E_{\pi(i),i}R$ is an eigenvector of \mathcal{H}_a associated to the eigenvalue $\sigma_{\pi(i)}$.
- For $1 \leq i < j \leq r$, $(E_{\pi(i),j} - E_{\pi(j),i})R$ is an eigenvector of \mathcal{H}_a associated to the eigenvalue 0.
- For $r+1 \leq k \leq \ell$, $1 \leq j \leq r$, $E_{\pi(k),j}R$ is an eigenvector of \mathcal{H}_a associated to the eigenvalue $\sigma_{\pi(k)} \left(1 - \frac{\sigma_{\pi(k)}}{\sigma_{\pi(j)}}\right)$.

Proof. Let $\pi \in \mathcal{P}$, $R \in \mathcal{O}_r$ and $\Delta \in \mathbb{R}^{\ell, r}$. Using the change of variables $M = SA$ and denoting \mathcal{H}_a and \mathcal{H}_m the Hessian of respectively f_a at $\Pi_\pi R$ and f_m at $S\Pi_\pi R$, we have the equality :

$$\mathcal{H}_a[\Delta, \Delta] = \mathcal{H}_m[S\Delta, S\Delta].$$

After normalizing the eigenvectors of \mathcal{H}_m given in Theorem 53, we obtain :

$$\begin{aligned} \mathcal{H}_a[\Delta R, \Delta R] &= \mathcal{H}_m[S\Delta R, S\Delta R] = \sum_{1 \leq i < j \leq r} \left\langle (\sigma_{\pi(i)}^{-1} + \sigma_{\pi(j)}^1)^{-\frac{1}{2}} S^{-1}(E_{\pi(i),j} + E_{\pi(j),i}), S\Delta \right\rangle^2 \\ &\quad + \sum_{1 \leq i \leq r} \left\langle \sigma_{\pi(i)}^{\frac{1}{2}} S^{-1} E_{\pi(i),i}, S\Delta \right\rangle^2 \\ &\quad + \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \left(1 - \frac{\sigma_{\pi(k)}}{\sigma_{\pi(j)}}\right) \langle E_{\pi(k),j}, S\Delta \rangle^2 \\ &= \sum_{1 \leq i < j \leq r} (\sigma_{\pi(i)}^{-1} + \sigma_{\pi(j)}^1)^{-1} \langle E_{\pi(i),j} + E_{\pi(j),i}, \Delta \rangle^2 \\ &\quad + \sum_{1 \leq i \leq r} \sigma_{\pi(i)} \langle E_{\pi(i),i}, \Delta \rangle^2 \\ &\quad + \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \sigma_{\pi(k)} \left(1 - \frac{\sigma_{\pi(k)}}{\sigma_{\pi(j)}}\right) \langle E_{\pi(k),j}, \Delta \rangle^2. \end{aligned}$$

□

As a direct corollary of Theorem 55, we have the following result.

Corollary 56. *With the notations used in Equation (68), an optimum $\tilde{I}R$ of f_a corresponds to the identity permutation $\pi = Id$. At an optimum, the largest eigenvalue of the Hessian \mathcal{H}_a is σ_1 and $\sigma_{\pi(\ell)} \left(1 - \frac{\sigma_{\pi(r+1)}}{\sigma_{\pi(r)}}\right) > 0$ is a lower bound of the positive eigenvalues of \mathcal{H}_a .*

The following result is also a straightforward corollary of Theorem 55.

Corollary 57. *All full-rank critical points that are not global minima are strict saddle points i.e. the Hessian at these points has a negative eigenvalue.*

Proof. Consider $R \in \mathcal{O}_r$ and a permutation $\pi : \llbracket 1; \ell \rrbracket \rightarrow \llbracket 1; \ell \rrbracket$ such that $\pi(1) < \dots < \pi(r)$ and simultaneously $\pi(r+1) < \dots < \pi(\ell)$ while $\pi \neq Id$. Necessarily, $\pi(r+1) < \pi_r$ and $\sigma_{\pi(r+1)} \left(1 - \frac{\sigma_{\pi(r+1)}}{\sigma_{\pi(r)}}\right) < 0$ is an eigenvalue of \mathcal{H}_a at $\Pi_\pi R$ by Theorem 55. \square

We can now prove Theorem 4.

Proof of Theorem 4. Consider a minimum $\tilde{I}R$ of f_a with $R \in \mathcal{O}_r$. From Lemma 3, we know that

$$\mathcal{C}_a(R) = \left\{ \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} R \mid A_1 \in \mathcal{S}_r^+, A_2 \in \mathbb{R}^{\ell-r, r} \right\}.$$

We denote the subspace spanned by $\mathcal{C}_a(R)$

$$\mathcal{E}_R^+ := \text{span} [\mathcal{C}_a(R)] = \left\{ \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} R \mid A_1 \in \mathcal{S}_r, A_2 \in \mathbb{R}^{\ell-r, r} \right\},$$

where \mathcal{S}_r is the set of symmetric matrices in $\mathbb{R}^{r, r}$. We know from Theorem 55 that \mathcal{E}_R^+ is exactly the subspace spanned by the eigenvectors of the Hessian $\mathcal{H}_{\tilde{I}R}$ of f_a at $\tilde{I}R$ associated to positive eigenvalues. Let $\sigma_{\min} := \sigma_\ell \left(1 - \frac{\sigma_{r+1}}{\sigma_r}\right)$. As pointed out in Corollary 56, σ_{\min} is a lower bound for the positive eigenvalue of the Hessian $\mathcal{H}_{\tilde{I}R}$. Thus, for all $M \in \text{span} (\mathcal{C}_a(R))$,

$$\text{Vec}(M)^T \mathcal{H}_{\tilde{I}R} \text{Vec}(M) \geq \sigma_{\min} \|M\|_F^2,$$

where $\text{Vec}(M) \in \mathbb{R}^{\ell, r}$ is the vectorization of $M \in \mathbb{R}^{\ell, r}$. Given the form of the Hessian for the trace norm in Proposition 6 of (Grave et al., 2011) that is recalled in Proposition 66, the existence of continuous bases for the singular subspaces (Stewart, 2012) of $S^2 \tilde{I}$ and the converse of Taylor's Theorem in (Oliver, 1954), we obtain that the Hessian of f_a is continuous at $\tilde{I}R$. Therefore, for any $\gamma < 1 < \delta$, there exists $\alpha > 0$ such that for all $M \in \mathcal{E}_R^+$ and $A \in \mathcal{B}(\tilde{I}R, \alpha) \cap \mathcal{E}_R^+$ where $\mathcal{B}(\tilde{I}R, \alpha)$ is the ball with center $\tilde{I}R$ and radius α , we have

$$\delta \sigma_1 \|M\|_F^2 \geq \text{Vec}(M)^T \mathcal{H}_A \text{Vec}(M) \geq \gamma \sigma_{\min} \|M\|_F^2. \quad (69)$$

Consider two elements $M, N \in \mathcal{B}(\tilde{I}R, \alpha)$. The Taylor expansions gives

$$\begin{aligned} f_a(N) &= f_a(M) + \langle \nabla f_a(M), N - M \rangle + \frac{1}{2} \int_0^1 \text{Vec}(N - M)^T \mathcal{H}_{tN + (1-t)M} \text{Vec}(N - M) dt \\ &\geq f_a(M) + \langle \nabla f_a(M), N - M \rangle + \frac{\gamma \sigma_{\min}}{2} \|N - M\|_F^2. \end{aligned}$$

This inequality implies that f_a is $\gamma \sigma_{\min}$ -strongly convex in $\mathcal{B}(\tilde{I}R, \alpha) \cap \mathcal{E}_R^+$. We conclude by defining a sublevel set \mathcal{V}_a inside $\cup_{R \in \mathcal{O}_r} \mathcal{B}(\tilde{I}R, \alpha)$.

Similarly, we could show from Equation (69) that for any $A, A' \in \mathcal{V}_a$ such that $[A, A'] \subset \mathcal{V}_a$, the function f_a has $\delta \sigma_1$ -Lipschitz gradients on $[A, A']$. Unfortunately, we can not deduce from this observation that f_a has $\delta \sigma_1$ -Lipschitz gradients or is $\delta \sigma_1$ -smooth in \mathcal{V}_a since the latter might be nonconvex. However, as in

Equation 24 of Lemma 34, s_1^2 being the largest eigenvalue of S^2 , we have for any $A, A' \in \mathbb{R}^{\ell, r}$, such that f_a is differentiable at A ,

$$f_a(A') \leq f_a(A) + \langle \nabla f_a(A), A' - A \rangle + \frac{s_1^2}{2} \|A' - A\|_F^2.$$

Therefore, the function f_a is s_1^2 -smooth. □

Remark 58. Note that the assumption $s_r > s_{r+1}$ is essential here. In order to highlight its importance, we can give an example to demonstrate that Theorem 4 would not be true if this assumption was not satisfied. Consider $Y = X = I_2 \in \mathbb{R}^{2,2}$ and $r = 1$. Here, the assumptions $s_r > s_{r+1}$ is violated since $s_1 = s_2 = 1$. The cones are $\mathbb{R}_+ \times \mathbb{R}$ and $\mathbb{R}_- \times \mathbb{R}$. The matrix $U = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ is an optimum of

$$\min_{U \in \mathbb{R}^{2,1}} \frac{1}{2} \|XU\|_F^2 - \|Y^T XU\|_* = \min_{U \in \mathbb{R}^{2,1}} \frac{1}{2} \|U\|_F^2 - \|U\|_*.$$

However, in the direction $\Delta_\alpha := \begin{bmatrix} 0 \\ \alpha \end{bmatrix}$, there is no strong convexity. Indeed we have

$$\frac{1}{2} \|X(U + \Delta_\alpha)\|_F^2 = \frac{1}{2} \|U + \Delta_\alpha\|_F^2 = \frac{1}{2}(1 + \alpha^2)$$

and

$$\|Y^T X(U + \Delta_\alpha)\|_* = \|U + \Delta_\alpha\|_* = \sqrt{1 + \alpha^2} = 1 + \frac{1}{2}\alpha^2 + o(\alpha^2).$$

By taking the difference of these two equations we prove that there is no second order dependance and consequently no strong convexity in the direction $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$. It could have been seen directly with Theorem 55 : with

$r = 1, \ell = 2, \pi = Id, E_{\pi(2),1} = E_{2,1} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ is an eigenvector associated to the eigenvalue $\sigma_1(1 - \frac{\sigma_2}{\sigma_1}) = 0$ since $\sigma_1 = \sigma_2 = 1$.

H.3 Proof of Corollary 5

Here, we do not assume that $X^T X$ is invertible and prove a more general result. We show that for any $R \in \mathcal{O}_r$ and $N \in \mathbb{R}^{p-m, r}$, the function f restricted to the affine cone $\mathcal{C}(R, N) = \tau(\mathcal{C}_a(R), \mathbb{R}^{m-\ell, r}, N)$, where τ is the function defined in Equation (51), is strongly convex in a neighborhood of the optimum $\tau(R, 0, N)$ of f . If we assumed that $X^T X$ is invertible, the proof would be very similar since we would have $m = p$ and the value of $f \circ \tau$ does not depend on N .

Given $R \in \mathcal{O}_r$ and $N \in \mathbb{R}^{p-m, r}$, consider U and U' in the same cone $\mathcal{C}(R, N)$ as $\tau(R, 0, N)$. Using the linear change of variables τ , we know that there exists $A, A' \in \mathcal{C}_a(R)$ and $C, C' \in \mathbb{R}^{m-\ell, r}$ such that :

$$\begin{aligned} U &= (X^T X)^{\frac{1}{2}} (PSA + P^\perp C) + K^\perp N \\ &= (X^T X)^{\frac{1}{2}} (PM + P^\perp C) + K^\perp N \text{ with } M = SA, \end{aligned}$$

and similarly $U' = (X^T X)^{\frac{1}{2}} (PM' + P^\perp C') + K^\perp N$ with $M' = SA'$.

We know from Equation (52) that:

$$f(U) = \frac{1}{2} \|M\|_F^2 - \|SM\|_* + \frac{1}{2} \|C\|_F^2.$$

In Theorem 53, we have computed the eigenvectors and the eigenvalues of $f_m : M'' \mapsto \frac{1}{2} \|M''\|_F^2 - \|SM''\|_*$ at $S\tilde{I}R$ which is a minimum of f_m . We invoke the same arguments as in the proof of Theorem 4 : given the

form of the Hessian for the trace norm in Proposition 6 of (Grave et al., 2011) that is recalled in Proposition 66, the existence of continuous bases for the singular subspaces (Stewart, 2012) of $S^2\tilde{I}$ and the converse of Taylor's Theorem in (Oliver, 1954), we obtain that the Hessian of f_m is continuous at $S\tilde{I}R$. Therefore, for any $\gamma < 1 < \delta$, there exists $\alpha > 0$ such that if $S^{-1}M, S^{-1}M' \in \mathcal{B}(\tilde{I}R, \alpha) \cap \mathcal{E}_R^+$ where $\mathcal{B}(\tilde{I}R, \alpha)$ is the ball with center $\tilde{I}R$ and radius α , we have

$$\frac{\gamma}{2}\left(1 - \frac{s_{r+1}^2}{s_r^2}\right)\|M' - M\|_F^2 \leq f_m(M) - f_m(M') - \langle \nabla f_m(M'), M - M' \rangle \leq \frac{\delta}{2}\|M' - M\|_F^2 \quad (70)$$

since the smallest positive eigenvalue of the Hessian of f_m at $S\tilde{I}R$ is $1 - \frac{s_{r+1}^2}{s_r^2}$ and the largest is 1.

The variables U and U' being obtained from (M, C, N) and (M', C', N) with a linear transformation, we can define a neighborhood $\mathcal{V}(R, N) \subset \mathcal{C}(R, N)$ of $\tau(\tilde{I}R, 0, N)$ such that $U, U' \in \mathcal{V}(R, N)$ if and only if $S^{-1}M, S^{-1}M' \in \mathcal{B}(\tilde{I}R, \alpha) \cap \mathcal{E}_R^+$ and then transfer Equation (70) to U and U' :

$$\begin{aligned} & \frac{\gamma}{2}\left(1 - \frac{s_{r+1}^2}{s_r^2}\right) \left[\|C - C'\|_F^2 + \frac{1}{2}\|M' - M\|_F^2 \right] \\ & \leq \frac{1}{2}\|C - C'\|_F^2 + \frac{\gamma}{2}\left(1 - \frac{s_{r+1}^2}{s_r^2}\right)\|M' - M\|_F^2 \\ & \leq f(U) - f(U') - \langle \nabla f(U'), U - U' \rangle. \end{aligned}$$

Also, since $U - U' = (X^T X)^{\frac{1}{2}}(P(M - M') + P^\perp(C - C'))$ we have the following inequality

$$\|U - U'\|_F^2 \leq d_{\max}^2 [\|M - M'\|_F^2 + \|C - C'\|_F^2],$$

where d_{\max} is the largest eigenvalue of $(X^T X)^{\frac{1}{2}}$. If $X^T X$ is invertible, $\frac{1}{d_{\max}^2}$ is the smallest eigenvalue of $X^T X$. Eventually, we obtain

$$\frac{\gamma}{2d_{\max}^2}\left(1 - \frac{s_{r+1}^2}{s_r^2}\right)\|U - U'\|_F^2 \leq f(U) - f(U') - \langle \nabla f(U'), U - U' \rangle.$$

Setting $\mu := \frac{\gamma}{2d_{\max}^2}\left(1 - \frac{s_{r+1}^2}{s_r^2}\right)$, we have proved that the restriction of f to the affine cone $\mathcal{C}(R, N)$ is μ -strongly convex in the neighborhood $\mathcal{V}(R, N)$ of the optimum $\tau(\tilde{I}R, 0, N)$. We conclude by defining a sublevel set $\mathcal{V}^0 \subset \cup_{R \in \mathcal{O}_r, N \in \mathbb{R}^{p-m, r}} \mathcal{V}(R, N)$ of the function f .

The L_X -smoothness of the function f is obtained directly from Equation (24) in Fact 34.

Remark 59. *Similarly, we can show from Equation (70) that there exists $M > L_X$ such that for any $U, U' \in \mathcal{V}^0$ with $[U, U'] \subset \mathcal{V}^0$, the function f has M -Lipschitz gradients on $[U, U']$, since the Hessian is bounded in \mathcal{V}^0 . Unfortunately, we can not deduce from this observation that f has M -Lipschitz gradients in \mathcal{V}^0 or is M -smooth in \mathcal{V}^0 since the latter might be nonconvex.*

H.4 Proof of Corollary 6

To extend to (SRRR) the result that we proved for (RRR), we assume that $X^T X$ is invertible.

Proof of Corollary 6. Let $\mu < \nu_X \left(1 - \frac{s_r^2}{s_{r+1}^2}\right)$ where ν_X is the smallest eigenvalue of $X^T X$ and \mathcal{V}^0 be defined as in Corollary 5. As $X^T X$ is invertible, we know from the orthogonal invariance of $f(U)$ and $\lambda\|U\|_{1,2}$ that for any $R \in \mathcal{O}_r$, a minimum of $F^\lambda(U) = f(U) + \lambda\|U\|_{1,2}$ is attained in each cone $\mathcal{C}(R)$. Theorem 6.4 of Bonnans and Shapiro (1998) guarantees, if its conditions are satisfied, the existence of $\check{\lambda}$ such that for any $R \in \mathcal{O}_r$, the minimum in each cone $\mathcal{C}(R)$ depends continuously on $\lambda \in [0, \check{\lambda})$. The assumptions of the Theorem 6.4 are indeed satisfied and we detail those below :

- (a) The objective F^λ of (SRRR) is locally strongly convex on the cone $\mathcal{C}(I_r)$ around the minimum : indeed, the restriction to $\mathcal{C}(I_r)$ of $f : U \mapsto \frac{1}{2}\|XU\|_F^2 - \|Y^T XU\|_*$ is strongly convex according to Corollary 5 and $\lambda\|U\|_{1,2}$ is convex.
- (b) For every fixed λ in some interval $[0, \tilde{\lambda})$, f is locally Lipschitz with a constant that does not depend on λ and the group-Lasso norm is Lipschitz.
- (c) The difference $F^\lambda - F^0 = \lambda\|\cdot\|_{1,2}$ is locally Lipschitz with a constant $\sqrt{p}\lambda$ which is $O(\lambda)$.

Thus, according to Theorem 6.4 of Bonnans and Shapiro (1998), there exists $0 < \bar{\lambda} < \tilde{\lambda}$ such that for any $0 \leq \lambda \leq \bar{\lambda}$, the optimum of (SRRR) in each cone remains in the neighborhood \mathcal{V}^0 where f is L_X -smooth and F^λ is μ -strongly convex, with the same constants as f for (RRR). To conclude and obtain Corollary 6, there only remains to define a new open sublevel set \mathcal{V}^λ of F^λ inside the sublevel set \mathcal{V}^0 of f . □

I Proofs for Section 5.3

I.1 Proof of Theorem 7

The sequence of inequalities to prove Theorem 7 is the same as in Proof B.1 of Csiba and Richtárik (2017) except for the line search condition that plays the role of their smoothness condition. Indeed, the result remains true if the function is not smooth as long as the condition (LS) is satisfied. Let $F^{\lambda,*}$ denote the minimum of F^λ . We define the optimality gap function

$$\xi : x \mapsto F^\lambda(x) - F^{\lambda,*}.$$

Given $t > 0$ and a point $x \in \mathbb{R}^d$, we have also defined

$$\begin{aligned} \tilde{f}_{t,x}(x') &:= f(x) + \langle \nabla f(x), x' - x \rangle + \frac{1}{2t} \|x' - x\|_F^2, \\ \tilde{F}_{t,x}^\lambda(x') &:= \tilde{f}_{t,x}(x') + \lambda h(x'), \end{aligned} \tag{71}$$

and x^+ is the unique minimum of the strongly convex function $\tilde{F}_{t,x}^\lambda$.

Proof of Theorem 7. Let $x \in \mathbb{R}^d$ and $t > 0$ such that the condition (LS) is satisfied i.e. $\tilde{F}_{t,x}^\lambda(x^+) \geq F^\lambda(x^+)$. We have

$$\xi(x^+) = F^\lambda(x^+) - F^{\lambda,*} \tag{72}$$

$$\leq \tilde{F}_{t,x}^\lambda(x^+) - F^{\lambda,*} \tag{73}$$

$$= f(x) + \lambda h(x) - F^{\lambda,*} + \langle \nabla f(x), x^+ - x \rangle + \frac{1}{2t} \|x^+ - x\|^2 + \lambda h(x^+) - \lambda h(x) \tag{74}$$

$$= \xi(x) + \min_{y \in \mathbb{R}^d} \left[\langle \nabla f(x), y - x \rangle + \frac{1}{2t} \|y - x\|^2 + \lambda h(y) - \lambda h(x) \right] \tag{75}$$

$$= \xi(x) - t\gamma_t(x) \tag{76}$$

$$= \xi(x) [1 - t\alpha_t(x)] \tag{77}$$

Equation (72) follows from the definition of ξ . We have Equation (73) since the condition (LS) is satisfied. Equation (74) comes from Equation (71). Equation (75) follows from the definition of x^+ , Equation (76) from the definition of γ_t and Equation (77) from the definitions of α_t and ξ . □

Remark 60. *A similar result would hold if we used stochastic block coordinate descent as in Lemma 13 of Csiba and Richtárik (2017), the proof would again follow Proof B.1 in Csiba and Richtárik (2017), with the same modification about the condition (LS).*

J Proofs for Section 5.4

J.1 Proof of Corollary 8

Let μ and \mathcal{V}^0 be defined as in Corollary 5. Let $R \in \mathcal{O}_r$ and $U \in \mathcal{C}(R) \cap \mathcal{V}^0$. According to Corollary 5, f is μ -strongly convex on $\mathcal{C}(R) \cap \mathcal{V}^0$. Since the minimal value f^* of f is attained on each cone, let $U^* \in \mathcal{C}(R)$ be an optimum of f . As $\mathcal{C}(R) \cap \mathcal{V}^0$ defines a sublevel set of the restriction of f to $\mathcal{C}(R)$ that is a convex function, it is a convex set. Therefore, the segment $[U^*, U]$ is included in $\mathcal{C}(R) \cap \mathcal{V}^0$.

As a μ -strongly convex function, the restriction $f|_{\mathcal{C}(R) \cap \mathcal{V}^0}$ of f to the convex set $\mathcal{C}(R) \cap \mathcal{V}^0$ satisfies

$$f|_{\mathcal{C}(R) \cap \mathcal{V}^0}(U^*) \geq f|_{\mathcal{C}(R) \cap \mathcal{V}^0}(U) + \langle \nabla f|_{\mathcal{C}(R) \cap \mathcal{V}^0}(U), U^* - U \rangle + \frac{\mu}{2} \|U^* - U\|_F^2.$$

Since

$$\langle \nabla f|_{[U, U^*]}(U), U' - U \rangle = \lim_{s \rightarrow 0^+} \frac{f(U + s(U' - U)) - f(U)}{s} = \langle \nabla f(U), U' - U \rangle,$$

we obtain

$$\begin{aligned} f(U) - f^* &\leq \langle \nabla f(U), U - U^* \rangle - \frac{\mu}{2} \|U - U^*\|_F^2 \\ &= \frac{\mu}{2} \left(\frac{1}{\mu} \|\nabla f(U)\|_F^2 - \|U - U^*\|_F^2 - \frac{1}{\mu} \|\nabla f(U)\|_F^2 \right) \\ &\leq \frac{1}{2\mu} \|\nabla f(U)\|_F^2. \end{aligned}$$

J.2 Proof of Corollary 9

First, we need to introduce the following lemma. It is a light modification of Theorem 15 of Csiba and Richtárik (2017). Apart from the substitution of the Lipschitz constant with $\frac{1}{t}$, the proof follows Proof B.2 of Csiba and Richtárik (2017).

Lemma 61. *Let $\lambda \geq 0$, $\mu \geq 0$, $\mathcal{C} \subset \mathbb{R}^{p,r}$ a convex set, $f : \mathbb{R}^{p,r} \rightarrow \mathbb{R}$ be a differentiable function such that its restriction to \mathcal{C} is μ -strongly convex, $h : \mathbb{R}^{p,r} \rightarrow \mathbb{R}$ be a convex function and $F^\lambda = f + \lambda h$. We denote \bar{f} , \bar{h} and \bar{F}^λ the restrictions of f , h and F^λ to \mathcal{C} . $F^{\lambda,*}$ denotes the optimal value of \bar{F}^λ in \mathcal{C} . Given $U, U' \in \mathcal{C}$ and $t > 0$, we denote*

$$\tilde{\bar{F}}^\lambda(U') := \bar{f}(U) + \langle \nabla \bar{f}(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \bar{h}(U') \quad (78)$$

$$\tilde{\gamma}_t(U) := -\frac{1}{t} \min_{U' \in \mathcal{C}} \left[\tilde{\bar{F}}^\lambda(U') - \bar{F}^\lambda(U) \right]. \quad (79)$$

Let $U \in \mathcal{C}$, $t > 0$ and $U_+ = \operatorname{argmin}_{U' \in \mathcal{C}} \left[\tilde{\bar{F}}^\lambda(U') - \bar{F}^\lambda(U) \right]$. We have

$$\tilde{\gamma}_t(U) \geq \min \left(\frac{1}{2t}, \mu \right) \left[\bar{F}^\lambda(U) - \bar{F}^{\lambda,*} \right].$$

Proof. Let $U \in \mathcal{C}$ such that $\bar{F}^\lambda > \bar{F}^{\lambda,*}$, $t > 0$ and $U_+ = \operatorname{argmin}_{U' \in \mathcal{C}} \left[\tilde{\bar{F}}^\lambda(U') - \bar{F}^\lambda(U) \right]$. We have

$$t\tilde{\gamma}_t(U) = -\min_{U' \in \mathcal{C}} \left[\langle \nabla \bar{f}(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \bar{h}(U') - \bar{h}(U) \right] \quad (80)$$

$$\begin{aligned} &= \bar{F}^\lambda(U) - \min_{U' \in \mathcal{C}} \left[\bar{f}(U) + \langle \nabla \bar{f}(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \bar{h}(U') \right] \\ &\geq \bar{F}^\lambda(U) - \min_{U' \in \mathcal{C}} \left[\bar{f}(U') - \frac{\mu}{2} \|U' - U\|_F^2 + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \bar{h}(U') \right] \end{aligned} \quad (81)$$

$$= \bar{F}^\lambda(U) - \min_{U' \in \mathcal{C}} \left[\bar{F}^\lambda(U') - \frac{1}{2} \left(\mu - \frac{1}{t} \right) \|U' - U\|_F^2 \right] \quad (82)$$

Equation (80) follows from Equation (78) and Equation (79). Equation (81) is due to the μ -strong convexity of \bar{f} . We denote $U^* \in \mathcal{C}$ the optimum of \bar{F}^λ and for all $U' \in \mathcal{C}$, $\xi(U') := \bar{F}^\lambda(U') - \bar{F}^{\lambda,*}$. Let $0 \leq \delta \leq 1$, setting $U' = U + \delta(U^* - U)$ in Equation (82), we obtain :

$$\begin{aligned} t\bar{\gamma}_t(U) &\geq \bar{F}^\lambda(U) - \bar{F}^\lambda(\delta U^* + (1-\delta)U) + \frac{1}{2}\delta^2 \left(\mu - \frac{1}{t} \right) \|U^* - U\|_F^2 \\ &\geq \bar{F}^\lambda(U) - \delta \bar{F}^\lambda(U^*) - (1-\delta)\bar{F}^\lambda(U) + \frac{1}{2} \left[\mu\delta(1-\delta) + \delta^2 \left(\mu - \frac{1}{t} \right) \right] \|U^* - U\|_F^2 \end{aligned} \quad (83)$$

$$= \delta \left(\xi(U) + \frac{\mu}{2} \|U^* - U\|_F^2 \right) - \frac{\delta^2}{2t} \|U^* - U\|_F^2 \quad (84)$$

Equation (83) comes from the μ -strong convexity of \bar{F}^λ . We impose

$$\delta = \min \left(1, \frac{\xi(U) + \frac{\mu}{2} \|U - U^*\|_F^2}{\frac{1}{t} \|U - U^*\|_F^2} \right). \quad (85)$$

Consider the two possible values for δ in Equation (85). First, if $\frac{1}{t} \|U - U^*\|_F^2 \leq \xi(U) + \frac{\mu}{2} \|U - U^*\|_F^2$ we have $\delta = 1$ and

$$\left(\mu - \frac{1}{t} \right) \|U - U^*\|_F^2 \geq \left(\frac{\mu}{2} - \frac{1}{t} \right) \|U - U^*\|_F^2 \geq -\xi(U). \quad (86)$$

Combining Equation (84) with Equation (86) in the case $\delta = 1$, we obtain

$$t\bar{\gamma}_t(U) \geq \xi(U) + \frac{1}{2} \left(\mu - \frac{1}{t} \right) \|U^* - U\|_F^2 \geq \frac{1}{2} \xi(U). \quad (87)$$

Secondly, if $\frac{1}{t} \|U - U^*\|_F^2 \geq \xi(U) + \frac{\mu}{2} \|U - U^*\|_F^2$, we obtain with Equation (84)

$$t\bar{\gamma}_t(U) \geq \frac{\left(\xi(U) + \frac{\mu}{2} \|U^* - U\|_F^2 \right)^2}{\frac{2}{t} \|U^* - U\|_F^2}. \quad (88)$$

Therefore, with Equation (87) and Equation (88), we have

$$\begin{aligned} \bar{\gamma}_t(U) &\geq \min \left(\frac{1}{2t} \xi(U), \frac{\left(\xi(U) + \frac{\mu}{2} \|U^* - U\|_F^2 \right)^2}{2 \|U^* - U\|_F^2} \right) \\ &\geq \min \left(\frac{1}{2t} \xi(U), \frac{2\xi(U)\mu \|U^* - U\|_F^2}{2 \|U^* - U\|_F^2} \right) \\ &\geq \min \left(\frac{1}{2t}, \mu \right) \xi(U) = \min \left(\frac{1}{2t}, \mu \right) [\bar{F}^\lambda(U) - \bar{F}^{\lambda,*}]. \end{aligned} \quad (89)$$

Equation (89) comes from the inequality of arithmetic and geometric means. \square

We can now prove Corollary 9. Let \mathcal{V}^λ be the sublevel set defined in Corollary 6. Let $R \in \mathcal{O}_r$ and $U \in \mathcal{C}(R) \cap \mathcal{V}^\lambda$. According to Corollary 6, F^λ is μ -strongly convex on $\mathcal{C}(R) \cap \mathcal{V}^\lambda$. Since the minimal value $F^{\lambda,*}$ is attained on each cone, let $U^* \in \mathcal{C}(R)$ be an optimum of $F^{\lambda,*}$. As $\mathcal{C}(R) \cap \mathcal{V}^\lambda$ defines a sublevel set of the restriction of F^λ to $\mathcal{C}(R)$ that is a convex function, it is a convex set. Therefore, the segment $[U^*, U]$ is included in $\mathcal{C}(R) \cap \mathcal{V}^\lambda$.

We define for any $U' \in [U, U^*]$ the surrogate $(\tilde{F}^\lambda|_{[U, U^*]})_{t,x}(U')$ of the restriction of F^λ to $[U, U^*]$ like in section 5.3 :

$$(\tilde{F}^\lambda|_{[U, U^*]})_{t,U}(U') = f(U) + \langle \nabla f|_{[U, U^*]}(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|_F^2 + \lambda \|U'\|_{1,2}.$$

From Lemma 61, we obtain the following inequality for any $U' \in [U, U^*]$ such that the condition (LS) is satisfied :

$$-\frac{1}{t} \min_{U' \in [U, U^*]} \left[(\tilde{F}^\lambda|_{[U, U^*]})_{t, U}(U') - F^\lambda|_{[U, U^*]}(U) \right] \geq \min\left(\frac{1}{2t}, \mu\right) [F^\lambda(U) - F^{\lambda, *}] . \quad (90)$$

Since

$$\langle \nabla f|_{[U, U^*]}(U), U' - U \rangle = \lim_{s \rightarrow 0^+} \frac{f(U + s(U' - U)) - f(U)}{s} = \langle \nabla f(U), U' - U \rangle,$$

Inequality (90) becomes :

$$-\frac{1}{t} \min_{U' \in [U, U^*]} \left\{ \tilde{F}_{t, U}^\lambda(U') - F^\lambda(U) \right\} \geq \min\left(\frac{1}{2t}, \mu\right) [F^\lambda(U) - F^{\lambda, *}] .$$

The minimum over the segment being lower bounded by the minimum over the whole space, we deduce that

$$\gamma_t(U) \geq \min\left(\frac{1}{2t}, \mu\right) [F^\lambda(U) - F^{\lambda, *}] .$$

J.3 Proof of Corollary 10

Let $\lambda \geq 0$ and \mathcal{V}^λ be a non-empty sublevel set of F^λ such that for all $U \in \mathcal{V}^\lambda$, F^λ satisfies the t -strong proximal-PL inequality, as in Corollary 9. Let $k \geq 0$, $t_{k-1} > \frac{\beta}{L_X}$ and $U^k \in \mathcal{V}^\lambda$. If U^{k+1} and t_k are generated as in Algorithm 1 from $U^k \in \mathcal{V}^\lambda$ and t_{k-1} such that the (LS) condition $F^\lambda(U_{k+1}) \leq \tilde{F}_{t_k, U^k}^\lambda(U_{k+1})$ is satisfied, then we know from Fact 35 that the inequality $t_k > \frac{\beta}{L_X}$ is satisfied.

Besides, since

$$F^\lambda(U^{k+1}) \leq \tilde{F}_{t_k, U^k}^\lambda(U^{k+1}) = \min_{U' \in \mathbb{R}^{p, r}} \tilde{F}_{t_k, U^k}^\lambda(U') \leq \tilde{F}_{t_k, U^k}^\lambda(U^k) = F^\lambda(U^k)$$

and \mathcal{V}^λ is a sublevel set, it is clear that $U^{k+1} \in \mathcal{V}^\lambda$.

To obtain Equation (10), we can apply Theorem 7 since F^λ satisfies the t_k -strong proximal-PL inequality by Corollary 9 with $\alpha(t_k) := \min(\frac{1}{2t_k}, \mu)$:

$$\begin{aligned} F^\lambda(U^{k+1}) - F^{\lambda, *} &\leq [1 - t_k \alpha(t_k)] [F^\lambda(U^k) - F^{\lambda, *}] \\ &\leq \left[1 - \min\left(\frac{1}{2}, \mu t_k\right) \right] [F^\lambda(U^k) - F^{\lambda, *}] \\ &\leq [1 - \rho] [F^\lambda(U^k) - F^{\lambda, *}] \end{aligned}$$

where $\rho = \min(\frac{1}{2}, \beta \frac{\mu}{L_X}) \leq \min(\frac{1}{2}, \mu t_k)$.

K Supplementary Results and Proofs

K.1 Proof of Lemma 44

First, we prove the following fact.

Fact 62. *If U is a local minimizer of F^λ , then denoting*

$$V_U \in \operatorname{argmax}_{V \in \mathbb{R}^{k, r}: V^T V = I_r} \langle V, Y^T X U \rangle, \quad (91)$$

the matrix $W := UV_U^T \in \mathbb{R}^{p, k}$ has to be a local minimizer of $F_w : W \mapsto \frac{1}{2} \|XW\|_F^2 - \langle Y, XW \rangle + \lambda \|W\|_{1,2}$ among matrices of $\mathbb{R}^{p, k}$ whose rank is smaller than r .

Proof. We prove Fact 62 by contradiction, supposing that $W := UV_U^T$ is not a local minimizer. Without loss of generality, we can assume since F^λ is invariant when its argument is multiplied on the right by an orthogonal matrix that the columns of U are orthogonal. Indeed, if the SVD of U is $R_1 \Sigma R_2^T$, we can consider instead $U' = R_1 \Sigma$ and modify V_U accordingly. With this assumption, the right singular vectors of $W := UV_U^T$ with V_U defined by Equation (91) are exactly the columns of V_U . Since we supposed that W is not a local minimizer, there exists a sequence of matrices $(W_k)_{k \geq 0}$ with rank smaller than r and with limit W such that for each $k \geq 0$, $F_w(W_k) < F_w(W)$. For each $k \geq 0$, let V_k be a matrix with r columns containing at least the right singular vectors of W_k such that $V_k^T V_k = I_r$. In particular, using the continuity of the singular spaces (Stewart, 2012, Theorem V.2.7), we can impose that the sequence $(V_k)_{k \geq 0}$ has limit V_U . The sequence $(U_k)_{k \geq 0}$ defined for each $k \geq 0$ by $U_k = W_k V_k$ has limit U . For $k \geq 0$, this would mean $W_k = U_k V_k^T$ and

$$f(U_k) + \lambda \|U_k\|_{1,2} = F^\lambda(U_k) \leq F_w(U_k V_k^T) < F_w(UV_U^T) = f(U) + \lambda \|U\|_{1,2}.$$

This would contradict the fact that U is a local minimizer. Therefore $W = UV_U^T$ must be a local minimizer of F_w . \square

Proof of Lemma 44. We assume that for any $S \subset \{1, \dots, p\}$ of cardinality at least r , the matrix $X_S^T Y$ is full-rank, where X_S is the matrix formed by keeping the columns of X indexed by S . We prove Lemma 44 by contradiction, assuming that U is a local minimum which has at least r non-zero rows and a rank strictly smaller than r . Again, we denote $V_U \in \operatorname{argmax}_{V^T V = I_r} \langle V, Y^T X U \rangle$ and consider $W := UV_U^T$. First, we write without loss of generality

$$W = \begin{bmatrix} W_S \\ 0 \end{bmatrix}, \text{ with } |S| \geq r \text{ and } W_S \in \mathbb{R}^{|S|,k} \text{ only has non-zero rows.}$$

Secondly, $\operatorname{rank}(W_S) < r$ since W_S is extracted from W whose rank is smaller than r . According to Fact 62, W is a local minimizer of F_w among matrices with rank smaller than r so for any vectors $u \in \mathbb{R}^p$, $v \in \mathbb{R}^k$, the function $t \mapsto \frac{1}{2} \|Y - X(W + tuv^T)\|_F^2 + \lambda \|W + tuv^T\|_{1,2}$ has a minimum at zero. The first-order condition is :

$$u^T X^T (Y - XW)v + \lambda \sum_i u_i z_i^T v = 0,$$

where $u_i \in \mathbb{R}$ and denoting $W_{i,:}$ the i -th row of W , $z_i^T = \frac{W_{i,:}}{\|W_{i,:}\|_2}$ if $W_{i,:}$ is different from zero and z_i has a norm smaller than 1 otherwise. If we impose $v \in \operatorname{Ker} W_S$, we get $Wv = 0$ and $z_i^T v = 0$ for $i \in S$. Therefore we have,

$$u^T X^T Y v + \lambda \sum_{i \notin S} u_i z_i^T v = 0. \quad (92)$$

Since Equation (92) holds in particular for any $u \in \mathbb{R}^p$ such that $u_i = 0$ when $i \notin S$, we necessarily have for any $v \in \operatorname{Ker} W_S$,

$$X_S^T Y v = 0.$$

In other words, we have $\operatorname{Ker} W_S \subset \operatorname{Ker} X_S^T Y$. This implies that $\dim(\operatorname{Ker} X_S^T Y) \geq \dim(\operatorname{Ker} W_S) > k - r$ since W_S has rank strictly smaller than r . Therefore $X_S^T Y \in \mathbb{R}^{|S|,k}$ has rank strictly smaller than r . This is in contradiction with the assumption in Lemma 44. \square

K.2 Proof of Lemma 45

Let U^* be a full-rank local minimum of $F^\lambda : U \mapsto \frac{1}{2} \|XU\|_F^2 - \|Y^T XU\|_* + \lambda \|U\|_{1,2}$. Without loss of generality, we denote S the support of the rows of U^* and we write

$$U^* = \begin{bmatrix} U_S \\ 0_{p-m,r} \end{bmatrix}$$

where m is the number of non-zero rows of U and $U_S \in \mathbb{R}^{m,r}$. We also denote

$$X = [X_S \quad X_{S^c}]$$

with $X_S \in \mathbb{R}^{n,m}$ and $X_{S^c} \in \mathbb{R}^{n,p-m}$. Let $V \in \operatorname{argmin}_{V \in \mathbb{R}^{k,r}: V^T V = I_r} \langle Y^T X U^*, V \rangle$ and $G^\lambda : U \mapsto \frac{1}{2} \|XU\|_F^2 - \langle Y^T X U, V \rangle + \lambda \|U\|_{1,2}$. By Fact 25, we have on the one hand $G^\lambda \geq F^\lambda$ and on the other hand $G^\lambda(U^*) = F^\lambda(U^*)$ so U^* is a local minimum of G^λ . The first order conditions restricted to the rows in the set S are

$$X_S^T X_S U_S - X_S^T Y V + \lambda Z_S = 0 \quad (93)$$

where $Z_S := D U_S \in \mathbb{R}^{|S|,r}$ with $D := \operatorname{diag}(\frac{1}{\|U_1^*\|_2}, \dots, \frac{1}{\|U_m^*\|_2}) \in \mathbb{R}^{|S|,|S|}$ and $\|U_1^*\|_2, \dots, \|U_m^*\|_2$ are the norms of the rows of U_S . In particular, Equation (93) implies

$$U_S^T [X_S^T X_S + D] U_S = U_S^T X_S^T Y V.$$

Since we assumed that $|S| \geq r$, the matrix $U_S^T [X_S^T X_S + D] U_S$ has rank r . Necessarily, $U_S^T X_S^T Y = U^{*T} X^T Y$ also has rank r .

K.3 Proof of Lemma 46

Lemma 44 and Lemma 45 combined with Assumption $\mathcal{H}2$ ensure that for any limit point $U \in \bar{\mathcal{U}}$, the matrix $Y^T X U$ is full-rank. Since the set of limit points $\bar{\mathcal{U}}$ is closed and bounded, there exist $\zeta > 0$ and $\delta > 0$ such that for all $U \in \mathbb{R}^{p,r}$, $\operatorname{dist}(U, \bar{\mathcal{U}}) \leq \delta$ implies that the eigenvalues of $Y^T X U$ are lower bounded by ζ , where $\operatorname{dist}(U, \bar{\mathcal{U}})$ is the Euclidean distance between U and the compact set $\bar{\mathcal{U}}$. We denote $\mathcal{K}^\delta := \{U \in \mathbb{R}^{p,r} | \operatorname{dist}(U, \bar{\mathcal{U}}) \leq \delta\}$ and $\mathcal{K}^{\frac{\delta}{2}} := \{U \in \mathbb{R}^{p,r} | \operatorname{dist}(U, \bar{\mathcal{U}}) \leq \frac{\delta}{2}\}$.

Proposition 6 of (Grave et al., 2011) that is recalled in Proposition 66, describes the Hessian of the trace-norm at full-rank matrices : since for any $U \in \mathcal{K}^\delta$, the eigenvalues of $Y^T X U$ are lower bounded by ζ , there exists $M > 0$ such that the Hessian of f is bounded on \mathcal{K}^δ by M . Therefore, for any $U, U' \in \mathcal{K}^\delta$ such that $[U, U'] \subset \mathcal{K}^\delta$, we have

$$\|\nabla f(U) - \nabla f(U')\|_F \leq M \|U - U'\|_F. \quad (94)$$

Fact 35 and Lemma 36 ensure that $\lim_{k \rightarrow +\infty} \|U_{k+1} - U_k\|_F = 0$ so there exists $k_1 \geq 0$ such that for any $k \geq k_1$, $U_k \in \mathcal{K}^{\frac{\delta}{2}}$ and $\|U_{k+1} - U_k\|_F \leq \frac{\delta}{2}$. The triangle inequality implies that $[U_k, U_{k+1}] \subset \mathcal{K}^\delta$. Consequently we have, by Equation (94), for all $k \geq k_1$:

$$\|\nabla f(U_k) - \nabla f(U_{k+1})\|_F \leq M \|U_k - U_{k+1}\|_F.$$

K.4 Proof of Lemma 49

Let $A \in \mathbb{R}^{\ell,r}$ be a rank deficient matrix and $R_1 D R_2^T$ be a singular value decomposition of the matrix $S^2 A$. Since $S^2 A$ is rank deficient, we can assume $R_1 \in \mathbb{R}^{\ell,r-1}$, $D \in \mathbb{R}^{r-1,r-1}$ and $R_2 \in \mathbb{R}^{r,r-1}$. Up to a multiplication on the right by an orthogonal matrix, we can assume, using the orthogonal invariance of f_a , that

$$S^2 A = R_1 D [I_{r-1} \quad 0_{r-1}], \text{ where } I_{r-1} \in \mathbb{R}^{r-1,r-1}, 0_{r-1} \in \mathbb{R}^{r-1}.$$

Let $e_r \in \mathbb{R}^r$ be the vector whose components are 0 except for the last one that is 1. Let $t \in \mathbb{R}$ and $\tilde{a} \in \mathbb{R}^\ell$ be a unit-norm vector such that $S^2 \tilde{a}$ is orthogonal to the columns of R_1 and therefore to the columns of $S^2 A$. We have $\|\tilde{a} e_r^T\|_F = 1$. On the one hand, we can separate the Frobenius norm of $S(A + t \tilde{a} e_r^T)$ as follows,

$$\frac{1}{2} \|S(A + t \tilde{a} e_r^T)\|_F^2 = \frac{1}{2} \|SA\|_F^2 + \frac{1}{2} t^2 \|S \tilde{a} e_r^T\|_F^2 = \frac{1}{2} \|SA\|_F^2 + \frac{1}{2} t^2 \|S \tilde{a}\|_F^2 = \frac{1}{2} \|SA\|_F^2 + o(t).$$

On the other hand, for any $t \neq 0$, a singular value decomposition of $S^2(A + t \tilde{a} e_r^T)$ is

$$S^2(A + t \tilde{a} e_r^T) = \begin{bmatrix} R_1 & \frac{S^2 \tilde{a}}{\|S^2 \tilde{a}\|_F} \end{bmatrix} \begin{bmatrix} D & 0 \\ 0 & |t| \|S^2 \tilde{a}\|_F \end{bmatrix} \begin{bmatrix} I_{r-1} & 0_{r-1} \\ 0_{r-1}^T & \frac{t}{|t|} \end{bmatrix}.$$

We can therefore easily compute the trace norm of $S^2(A + t\tilde{a}e_r^T)$,

$$\|S^2(A + t\tilde{a}e_r^T)\|_* = \|S^2A\|_* + |t|\|S^2\tilde{a}\|_F \geq \|S^2A\|_* + |t|s_\ell^2,$$

where s_ℓ is the smallest eigenvalue of S . So finally, we obtain

$$f_a(A + t\tilde{a}e_r^T) \leq f_a(A) - s_\ell^2|t| + o(t).$$

K.5 Proof of Lemma 50

Proof. Let A be a critical point of $f_a : A \mapsto \frac{1}{2}\|SA\|_F^2 - \|S^2A\|_*$ and denote $V_A \in \operatorname{argmax}_{V \in \mathbb{R}^{\ell, r}: V^T V = I_r} \langle S, SAV^T \rangle$. We know from Lemma 49 that A is full-rank and applying Danskin's Theorem (Danskin, 1967), we have

$$\nabla f_a(A) = S^2(A - V_A) = 0. \quad (95)$$

Besides, writing $\Pi\Sigma R$ the singular value decomposition of S^2A with $\Pi \in \mathbb{R}^{\ell, r}$ a matrix whose columns are orthogonal, $\Sigma \in \mathbb{R}^{r, r}$ a diagonal matrix whose entries are denoted $\sigma_1, \dots, \sigma_\ell$ and $R \in \mathbb{R}^{r, r}$ an orthogonal matrix, we know that $A = V_A$ from Equation (95) and that $V_A = \Pi R$ by Fact 27. Therefore, we have

$$\begin{aligned} S^2A &= S^2\Pi R \\ \Rightarrow \Pi\Sigma R &= S^2\Pi R \quad \text{since } \Pi\Sigma R \text{ is the SVD of } S^2A, \\ \Rightarrow \Pi\Sigma &= S^2\Pi \quad \text{since } RR^T = I_r. \end{aligned} \quad (96)$$

Let $i \in \llbracket 1, r \rrbracket$, $w := (w_1, \dots, w_\ell)^T$ be the i -th column of Π . Equation (96) implies that

$$\begin{aligned} \sigma_i w &= \begin{bmatrix} s_1^2 w_1 \\ \vdots \\ s_\ell^2 w_\ell \end{bmatrix}, \\ \Rightarrow \forall j \in \llbracket 1, r \rrbracket, & (\sigma_i - s_j^2)w_j = 0. \end{aligned}$$

Since we assumed that s_1, \dots, s_ℓ are all different, only one w_j can be different from zero and must be 1 since w has norm 1. Given that the columns of the matrix Π are orthogonal and contain only one nonzero coefficient, up to a permutation of its columns, the matrix Π has the form given in Lemma 50.

With Lemma 47, we know that A is not a local maximum of f_a . If $A = \tilde{I}R$, we have proved in Lemma 1 that A is a global minimum. Now assume that $A = \Pi_\pi R$, with π and Π_π as in Equation (64) and that there exists $i \in \{1, \dots, r\}$ such that $\pi(i) > i$ and for all $i' < i$, $\pi(i') = i'$. We have

$$\frac{1}{2}\|S(A + te_i e_i^T R)\|_F^2 = \frac{1}{2}\|SA\|_F^2 + \frac{t^2}{2}s_i^2. \quad (97)$$

Since the Frobenius norm of the i -th column of $S^2(A + te_i e_i^T R)$ is $\sqrt{s_{\pi(i)}^4 + t^2 s_i^4}$ and the columns of $S^2(A + te_i e_i^T R)$ have disjoint supports, we also have

$$\|S^2(A + te_i e_i^T R)\|_* = \|S^2A\|_* - s_{\pi(i)}^2 + \sqrt{s_{\pi(i)}^4 + t^2 s_i^4} = \|S^2A\|_* - s_{\pi(i)}^2 + s_{\pi(i)}^2 \left(1 + \frac{t^2}{2} \frac{s_i^4}{s_{\pi(i)}^4}\right) + O(t^4) \quad (98)$$

Combining Equation (97) with Equation (98), we obtain

$$f_a(A + te_i e_i^T R) - f_a(A) = \frac{t^2 s_i^2}{2} \left(1 - \frac{s_i^2}{s_{\pi(i)}^2}\right) + O(t^4).$$

Since we have assumed that $\pi(i) > i$ and the eigenvalues of the matrix S are strictly decreasing, we have $\left(1 - \frac{s_i^2}{s_{\pi(i)}^2}\right) < 0$ and A is not a local minimum. □

K.6 Proof of Theorem 53

As in section G.3, we consider a permutation $\pi : \llbracket 1; \ell \rrbracket \rightarrow \llbracket 1; \ell \rrbracket$ such that simultaneously $\pi(1) < \dots < \pi(r)$ and $\pi(r+1) < \dots < \pi(\ell)$. We denote

$$\Pi_\pi := (1_{i=\pi(j)})_{1 \leq i \leq \ell, 1 \leq j \leq r} \in \mathbb{R}^{\ell, r},$$

and define for $i_0 \in \llbracket 1, \ell \rrbracket$ and $j_0 \in \llbracket 1, r \rrbracket$,

$$E_{i_0, j_0} = (1_{i=i_0, j=j_0})_{1 \leq i \leq \ell, 1 \leq j \leq r} \in \mathbb{R}^{\ell, r}.$$

We want to compute the Hessian \mathcal{H}_m of $f_m : M \mapsto \frac{1}{2} \|M\|_F^2 - \|SM\|_*$ at the matrix $M = S\Pi_\pi R$. It is well defined according to Proposition 6 in (Grave et al., 2011) since $SM = S^2\Pi_\pi R$ is full-rank. We recall this result below in Proposition 66. In order to introduce the different eigenvectors of the Hessian of f_m , we need the singular value decomposition and the polar decomposition of SM . Since $M = S\Pi_\pi R$ and $S^2\Pi_\pi = \Pi_\pi \text{diag}(s_{\pi(1)}^2, \dots, s_{\pi(r)}^2)$, a singular value decomposition of SM is given by

$$SM = \Pi_\pi \text{diag}(s_{\pi(1)}^2, \dots, s_{\pi(r)}^2) R, \quad \Pi_\pi^T \Pi_\pi = I_r \quad \text{and} \quad R^T R = I_r.$$

We have $s_{\pi(1)}^2 > \dots > s_{\pi(r)}^2$ because we assumed $s_1 > \dots > s_\ell > 0$ and $\pi(1) < \dots < \pi(r)$. Defining $V = \Pi_\pi R \in \mathbb{R}^{\ell, r}$ and $K = R^T \text{diag}(s_{\pi(1)}^2, \dots, s_{\pi(r)}^2) R \in \mathbb{R}^{r, r}$, we obtain the polar decomposition of SM :

$$SM = VK, \quad V^T V = I_r \quad \text{and} \quad K \in \mathcal{S}_r^{++} \quad (99)$$

with \mathcal{S}_r^{++} the set of positive-definite matrices in $\mathbb{R}^{r, r}$. We also denote $\mathcal{S}_r = \{H \in \mathbb{R}^{r, r} \mid H^T = H\}$ the set of symmetric matrices in $\mathbb{R}^{r, r}$.

First we focus on a set of directions where the restriction of f_m is exactly a quadratic strongly convex function.

Fact 63. *The restriction of $M' \mapsto \|SM'\|_*$ to the affine subspace $\{M + S^{-2}MH \mid H \in \mathcal{S}_r\}$ is linear in a neighborhood of M , its Hessian at M is zero. Consequently, the Hessian of $f_m : M \mapsto \frac{1}{2} \|M\|_F^2 - \|SM\|_*$ restricted to the subspace $T_{\mathcal{K}} := \{S^{-2}MH \mid H \in \mathcal{S}_r\}$ is exactly the identity. A basis for $T_{\mathcal{K}}$ is the concatenation of $(S^{-1}(E_{\pi(i), j} + E_{\pi(j), i})R)_{1 \leq i < j \leq r}$ with $(S^{-1}E_{\pi(i), i}R)_{1 \leq i \leq r}$.*

Proof. For any matrix \tilde{M} such that the polar decomposition of $S\tilde{M}$ has the form VB with $B \in \mathcal{S}_r^+$, we have $\|S\tilde{M}\|_* = \langle S\tilde{M}, V \rangle$. Indeed, if QDQ^T is a singular value decomposition of B with $Q \in \mathbb{R}^{r, r}$, $Q^T Q = I_r$ and $D \in \mathbb{R}^{r, r}$ a diagonal matrix, then $(VQ)DQ^T$ is a singular value decomposition of VB . Using Fact 25 and Fact 27, we have

$$\|S\tilde{M}\|_* = \langle S\tilde{M}, (VQ)Q^T \rangle = \langle S\tilde{M}, V \rangle.$$

Consequently, we have

$$f_m(\tilde{M}) = \frac{1}{2} \|\tilde{M}\|_F^2 - \langle S\tilde{M}, V \rangle.$$

In particular, for any $\Delta = S^{-1}VH$ with $H \in \mathcal{S}_r$ such that $K + H \in \mathcal{S}_r^+$, we have $M + \Delta = S^{-1}V(K + H)$ since $M = S^{-1}VK$ according to Equation (99) and

$$f_m(M + \Delta) = \frac{1}{2} \|M + \Delta\|_F^2 - \langle S(M + \Delta), V \rangle.$$

Therefore, the Hessian of $\Delta \mapsto f_m(M + \Delta)$ restricted to the subspace $T_{\mathcal{K}} := \{S^{-1}VH, H \in \mathcal{S}_r\}$ is locally the identity. Note that $S^{-1}V = S^{-2}S\Pi_\pi R = S^{-2}M$ since $V = \Pi_\pi R$ and $M = S\Pi_\pi R$ so

$$T_{\mathcal{K}} = \{S^{-2}MH, H \in \mathcal{S}_r\}.$$

We can also use $M = S\Pi_\pi R$ to write

$$\begin{aligned} T_{\mathcal{K}} &= \{S^{-1}\Pi_\pi RH, H \in \mathcal{S}_r\} \\ &= \{S^{-1}\Pi_\pi HR, H \in \mathcal{S}_r\}. \end{aligned} \quad (100)$$

For Equation (100), we have used the fact that for any orthogonal matrix $R \in \mathbb{R}^{r,r}$, the application $H \mapsto R^T HR$ is an automorphism of \mathcal{S}_r . We then obtain a basis for $T_{\mathcal{K}}$ using the fact that the concatenation of $(E_{i,j} + E_{j,i})_{1 \leq i < j \leq r}$ with $(E_{i,i})_{1 \leq i \leq r}$ is a basis of \mathcal{S}_r and for any $1 \leq i, j \leq r$, $\Pi_\pi E_{i,j} = E_{\pi(i),j}$. \square

Secondly, the invariance of f_m when its argument is multiplied on the right by an orthogonal matrix gives a set of directions included in the kernel of the Hessian.

Fact 64. *The subspace $T_{\mathcal{R}} := \{MT \mid T^T = -T, T \in \mathbb{R}^{r,r}\}$ is included in the Kernel of the Hessian of f_m at $M = S\Pi_\pi R$. Additionally, $T_{\mathcal{K}} \oplus^\perp T_{\mathcal{R}} = \{MF \mid F \in \mathbb{R}^{r,r}\}$ and a basis for $T_{\mathcal{R}}$ is $(S(E_{\pi(i),j} - E_{\pi(j),i})R)_{1 \leq i < j \leq r}$.*

Proof. Since M is a critical point of f_m which is invariant when its argument is multiplied on the right by an orthogonal matrix, then by (Li et al., 2016, Theorem 2), the subspace that is tangent to the manifold $\{MR' \mid R' \in \mathcal{O}_r\}$ is included in the Kernel of the Hessian of f_m at M . In Example 4, Li et al. (2016) show that this subspace is exactly $T_{\mathcal{R}} := \{MT \mid T \in \mathbb{R}^{r,r}, T^T = -T\}$. Since $M = S\Pi_\pi R$ and the set of antisymmetric matrices of $\mathbb{R}^{r,r}$ can be written $\{R^T TR \mid T \in \mathbb{R}^{r,r}, T^T = -T\}$, a basis for $T_{\mathcal{R}}$ is $(S(E_{\pi(i),j} - E_{\pi(j),i})R)_{1 \leq i < j \leq r}$.

To show that $\{MF \mid F \in \mathbb{R}^{r,r}\}$ can be decomposed with the given orthogonal sum, it is first important to notice that

$$\begin{aligned} T_{\mathcal{K}} &= \{S^{-1}VH \mid H \in \mathcal{S}_r\} \\ &= \{MK^{-1}H \mid H \in \mathcal{S}_r\}. \end{aligned} \quad (101)$$

We have used Equation (99) to obtain Equation (101). It is then sufficient to notice that both $T_{\mathcal{K}} = \{MK^{-1}H \mid H \in \mathcal{S}_r\}$ and $T_{\mathcal{R}} = \{MT \mid T^T = -T, T \in \mathbb{R}^{r,r}\}$ are included in $\{MF, F \in \mathbb{R}^{r,r}\}$, they are also orthogonal given the bases that we have introduced and finally, their dimensions are respectively $\frac{r(r+1)}{2}$ and $\frac{r(r-1)}{2}$ since M is full-rank so their sum must be equal to $\{MF \mid F \in \mathbb{R}^{r,r}\}$ which is of dimension r^2 . \square

What remains to study is the eigenvectors and the corresponding eigenvalues of the Hessian of f_m at M in the subspace that is orthogonal to $\{MF \mid F \in \mathbb{R}^{r,r}\}$.

Fact 65. *For $r+1 \leq k \leq \ell$ and $1 \leq j \leq r$, the matrix $E_{\pi(k),j}R$ is an eigenvector of the Hessian of f_m restricted to the subspace $T_{V^\perp} := \{C \in \mathbb{R}^{\ell,r} \mid M^T C = 0\}$ and the corresponding eigenvalue is $1 - \frac{s_\pi^2(k)}{s_\pi^2(j)}$.*

To prove Fact 65, we use the following result.

Proposition 66. *(Grave et al., 2011, Proposition 6) Let $\ell \geq r$, $N \in \mathbb{R}^{\ell,r}$ be a full-rank matrix and $W\Sigma Z^T \in \mathbb{R}^{\ell,r}$ be its singular value decomposition, with $W \in \mathbb{R}^{\ell,r}$, $W^T W = I_r$, $\Sigma = \text{diag}(\sigma_1 \geq \dots \geq \sigma_r) \in \mathbb{R}^{r,r}$ with $\sigma_r > 0$, $Z \in \mathbb{R}^{r,r}$ and $Z^T Z = I_r$. Let $W_0 \in \mathbb{R}^{\ell,\ell-r}$ such that $W_0^T W_0 = I_{\ell-r}$ and $W^T W_0 = 0$. We denote $(w_i)_{1 \leq i \leq r}$ the columns of W , $(z_j)_{1 \leq j \leq r}$ the columns of Z and $(w_k)_{r+1 \leq k \leq \ell}$ the columns of W_0 . For any $\Delta \in \mathbb{R}^{\ell,r}$, we have :*

$$\begin{aligned} \|N + \Delta\|_* &= \|N\|_* + \langle WZ^T, \Delta \rangle \\ &+ \frac{1}{2} \sum_{1 \leq i \leq r, 1 \leq j \leq r} \frac{(w_i^T \Delta z_j - w_j^T \Delta z_i)^2}{2(\sigma_i + \sigma_j)} \\ &+ \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \frac{(w_k^T \Delta z_j)^2}{\sigma_j} + o(\|\Delta\|_F^2). \end{aligned} \quad (102)$$

Proof of Fact 65. Given a perturbation ΔR of the matrix M , we have

$$\|SM + S\Delta R\|_* = \|S^2\Pi_\pi R + S\Delta R\|_*, \quad (103)$$

$$= \|S^2\Pi_\pi + S\Delta\|_* \quad (104)$$

Equation (103) comes from $M = S\Pi_\pi R$ and we have Equation (104) since the trace-norm is orthogonal invariant. Thus, we apply Proposition 66 for a perturbation $S\Delta$ of the matrix $S^2\Pi_\pi$ whose singular value decomposition is $\Pi_\pi \text{diag}(s_{\pi(1)}^2 > \dots > s_{\pi(r)}^2)$. With the notations of Proposition 66, this corresponds to $W = \Pi_\pi$, $\Sigma = \text{diag}(s_{\pi(1)}^2 > \dots > s_{\pi(r)}^2)$ and $Z = I_r$. Let $W_0 \in \mathbb{R}^{\ell, \ell-r}$ be the matrix whose columns w_k are the $e_{\pi(k)}$ for $r+1 \leq k \leq \ell$, then $W_0^T W_0 = I_{\ell-r}$ and $W^T W_0 = 0$. We have

$$\begin{aligned} \|S^2\Pi_\pi + S\Delta\|_* &= \|S^2\Pi_\pi\|_* + \langle WZ^T, S\Delta \rangle \\ &+ \frac{1}{2} \sum_{1 \leq i \leq r, 1 \leq j \leq r} \frac{(w_i^T S\Delta z_j - w_j^T S\Delta z_i)^2}{2(s_{\pi(i)}^2 + s_{\pi(j)}^2)} \\ &+ \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \frac{(w_k^T S\Delta z_j)^2}{s_{\pi(j)}^2} + o(\|\Delta\|_F^2) \\ &= \|S^2\Pi_\pi\|_* + \langle \Pi_\pi, S\Delta \rangle \\ &+ \frac{1}{2} \sum_{1 \leq i \leq r, 1 \leq j \leq r} \frac{(s_{\pi(i)} e_{\pi(i)}^T \Delta e_j - s_{\pi(j)} e_{\pi(j)}^T \Delta e_i)^2}{2(s_{\pi(i)}^2 + s_{\pi(j)}^2)} \\ &+ \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2} (e_{\pi(k)}^T \Delta z_j)^2 + o(\|\Delta\|_F^2). \end{aligned}$$

Note that in the first sum, Δ only intervenes through a product with the transpose of an element $e_{\pi(i)}$ that belongs to $\text{Im } M$. Since we already studied the effect of the Hessian on the subspace $\{MF \mid F \in \mathbb{R}^{r,r}\}$ in Fact 63 and Fact 64, we focus on the effect of the Hessian in the orthogonal subspace that is described in the second sum. Given a perturbation Δ of the form $W_0 F Z^T$ with $F \in \mathbb{R}^{\ell-r, r}$, we have on the one hand

$$\begin{aligned} \|S^2\Pi_\pi + S\Delta\|_* &= \|S^2\Pi_\pi\|_* + \langle \Pi_\pi, S\Delta \rangle \\ &+ \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2} (e_{\pi(k)}^T W_0 F Z^T z_j)^2 + o(\|\Delta\|_F^2) \\ &= \|S^2\Pi_\pi\|_* + \langle \Pi_\pi, S\Delta \rangle \\ &+ \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2} F_{k-r, j}^2 + o(\|\Delta\|_F^2). \end{aligned} \quad (105)$$

Equation (105) comes from $e_{\pi(k)}^T W_0 = (1_{i=k-r})_{1 \leq i \leq \ell-r}^T$ and $Z^T z_j = (1_{i=j})_{1 \leq i \leq r}$.

On the other hand,

$$\begin{aligned} \frac{1}{2} \|M + \Delta R\|_F^2 &= \frac{1}{2} \|M\|_F^2 + \frac{1}{2} \|\Delta R\|_F^2 + \langle M, \Delta R \rangle \\ &= \frac{1}{2} \|M\|_F^2 + \frac{1}{2} \|\Delta\|_F^2 + \langle S\Pi_\pi R, \Delta R \rangle \end{aligned} \quad (106)$$

$$= \frac{1}{2} \|M\|_F^2 + \frac{1}{2} \|W_0 F Z^T\|_F^2 + \langle S\Pi_\pi, \Delta \rangle \quad (107)$$

$$= \frac{1}{2} \|M\|_F^2 + \frac{1}{2} \sum_{1 \leq i \leq \ell-r, 1 \leq j \leq r} F_{i,j}^2 + \langle S\Pi_\pi, \Delta \rangle. \quad (108)$$

Equation (106) follows from $M = S\Pi_\pi R$, Equation (107) from $\Delta = W_0 F Z^T R$ and Equation (108) from $W_0^T W_0 = I_{\ell-r}$ and $Z = I_r$. Combining Equation (105) with Equation (108), we obtain for $\Delta = W_0 F Z^T$

$$f_m(M + \Delta R) = f_m(M) + \frac{1}{2} \sum_{r+1 \leq k \leq \ell, 1 \leq j \leq r} \left(1 - \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2}\right) F_{k-r,j}^2 + o(\|\Delta\|_F^2).$$

Since $W_0 \in \mathbb{R}^{\ell, \ell-r}$ is the matrix whose columns are the $e_{\pi(k)}$ for $r+1 \leq k \leq \ell$ and $Z = I_r$, we obtain the last eigenvectors of the Hessian of f_m : for $r+1 \leq k \leq \ell$ and $1 \leq j \leq r$, the matrix $E_{\pi(k),j} R$ is an eigenvector associated to the eigenvalue $1 - \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2}$. \square

Remark 67. Note that we could have directly used Equation (102) to prove simultaneously Fact 63, Fact 64 and Fact 65 but we believe that the proposed analysis helps understanding the structure of the eigenspaces.

Eventually, we have proved that the Hessian of f_m at M is block diagonal on the three orthogonal subspaces :

- $T_{\mathcal{K}} := \{S^{-2}MH \mid H \in \mathcal{S}_r\}$ where the eigenvalues are all equal to 1.
- $T_{\mathcal{R}} := \{MT \mid T^T = -T\}$ where the eigenvalues are all 0.
- $T_{V^\perp} := \{W_0 C \mid C \in \mathbb{R}^{\ell-r,r}\}$ where the eigenvalues are the $1 - \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2}$ for $r+1 \leq k \leq \ell$, $1 \leq j \leq r$.

We summarize the eigenvectors of the Hessian of $f_m : M' \mapsto \frac{1}{2}\|M'\|_F^2 - \|SM'\|_*$ at $M = S\Pi_\pi R$ in the table below.

Eigenvectors and Eigenvalues of the Hessian of $f_m : M' \mapsto \frac{1}{2}\ M'\ _F^2 - \ SM'\ _*$ at $M = \Pi_\pi R$			
Indices	Number of elements	Eigenvectors	Eigenvalues
$1 \leq i \leq r$	r	$S^{-1}E_{\pi(i),i}R$	1
$1 \leq i < j \leq r$	$\frac{r(r-1)}{2}$	$S^{-1}(E_{\pi(i),j} + E_{\pi(j),i})R$	1
$1 \leq i < j \leq r$	$\frac{r(r-1)}{2}$	$S(E_{\pi(i),j} - E_{\pi(j),i})R$	0
$r+1 \leq k \leq \ell, 1 \leq j \leq r$	$r(\ell-r)$	$E_{\pi(k),j}R$	$1 - \frac{s_{\pi(k)}^2}{s_{\pi(j)}^2}$

L KL with exponent $\frac{1}{2}$

As announced at the end of Section 5.4, we show in Section L.1 that the geometric structure leveraged in Corollary 9 can be used to prove that F^λ has the KL property with exponent $\frac{1}{2}$ near the set of optima. While in the core of the article, we proposed a direct proof of Corollary 10 based on Corollary 9 and Theorem 7, we present in Section L.2 an application of the framework developed by Csiba and Richtárik (2017) to show that the KL property with exponent $\frac{1}{2}$ (instead of the PL inequality) also leads to linear convergence. The proofs

appear simpler than the ones encountered in (Attouch and Bolte, 2009; Attouch et al., 2013; Chouzenoux et al., 2014; Frankel et al., 2015) as the algorithms considered in these papers are more general while we restrain our study to the proximal gradient algorithm with line search.

L.1 KL- $\frac{1}{2}$ on cones for (RRR / SRRR)

We assume that $X^T X$ is invertible. Let $\text{span } \mathcal{C}(I_r)$ be the subspace spanned by $\mathcal{C}(I_r) = \tau(\mathcal{C}_a(I_r), \mathbb{R}^{p-\ell, r})$ with τ defined in Equation (6) and $\mathcal{C}_a(I_r)$ defined in Equation (8). Let $F_{I_r}^\lambda$ be the restriction of F^λ to $\mathcal{C}(I_r)$: it is defined for any $U \in \text{span } \mathcal{C}(I_r)$ as $F_{I_r}^\lambda(U) = F^\lambda(U)$ if $U \in \mathcal{C}(I_r)$ and $F_{I_r}^\lambda(U) = +\infty$ otherwise. From the structure described in Corollary 6, and since τ is a linear invertible change of variables, we know that $F_{I_r}^\lambda \circ \tau$ is strongly convex in a neighborhood of $(\tilde{I}, 0_{p-\ell, r}) = \left(\begin{bmatrix} I_r \\ 0_{\ell-r, r} \end{bmatrix}, 0_{p-\ell, r} \right)$ included in $\mathcal{C}_a(I_r)$.

Fact 68. *Let F be a proper lower semi-continuous function. If F is μ -strongly convex in a set $\mathcal{V} \subset \mathbb{R}^d$ then given $x^* \in \mathcal{V}$, F has the Kurdyka-Lojasiewicz property at $x^* \in \text{dom } \partial F$ with exponent $1/2$: there exists $\eta > 0$ and a neighborhood \mathcal{U} of x^* such that for all $x \in \mathcal{U} \cap \{y \mid F(x^*) < F(y) < F(x^*) + \eta\}$, we have*

$$\frac{c}{\sqrt{F(x) - F(x^*)}} \text{dist}(0, \partial F(x)) \geq 1. \quad (109)$$

Proof. Let $x^* \in \mathcal{V}$. First, if $0 \notin \partial F(x^*)$, then by Lemma 2 of Attouch et al. (2010), there is $c > 0$ and a neighborhood \mathcal{U} of x^* such that for any $x \in \mathcal{U}$, we have

$$\text{dist}(0, \partial F(x)) \geq \frac{1}{c} \quad \text{and} \quad F(x) - F(x^*) \leq 1,$$

so Equation (109) holds for any $x \in \mathcal{U}$.

Secondly, assume that $0 \in \partial F(x^*)$. Let $x \in \mathcal{V}$ such that $F(x) > F(x^*)$ and $v \in \partial F(x)$. Since F is μ -strongly convex, we have:

$$\begin{aligned} F(x) - F(x^*) &\leq \langle v, x - x^* \rangle - \frac{\mu}{2} \|x - x^*\|^2 \\ &= \frac{\mu}{2} \left[\frac{1}{\mu^2} \|v\|^2 - \frac{1}{\mu^2} \|v\|^2 + 2 \left\langle \frac{1}{\mu} v, x - x^* \right\rangle - \|x - x^*\|^2 \right] \\ &= \frac{\mu}{2} \left[\frac{1}{\mu^2} \|v\|^2 - \|x - x^* - \frac{1}{\mu} v\|^2 \right] \\ &\leq \frac{1}{2\mu} \|v\|^2. \end{aligned}$$

Therefore, we obtain Equation (109) with $c = \frac{1}{\sqrt{2\mu}}$:

$$\frac{1}{\sqrt{2\mu}} \frac{1}{\sqrt{F(x) - F(x^*)}} \text{dist}(0, \partial F(x)) \geq 1.$$

□

Since $F_{I_r}^\lambda$ is strongly convex, it is a KL- $\frac{1}{2}$ function by Fact 68. This is key to apply the following result.

Theorem 69. *(Theorem 3.2 in Li and Pong, 2017) Consider $a \geq b \geq 1$, $g: \mathbb{R}^b \rightarrow \mathbb{R}$ a proper closed function and $h: \mathbb{R}^a \rightarrow \mathbb{R}^b$ a continuously differentiable mapping. Suppose in addition that g is a KL function with exponent $\alpha \in [0, 1)$ and the Jacobian $Jh(\bar{x}) \in \mathbb{R}^{b,a}$ is a surjective mapping at some $\bar{x} \in \text{dom } g \circ h$. Then $g \circ h$ has the KL property at \bar{x} with exponent α .*

Let $I_{p-\ell,r} : \mathbb{R}^{p-\ell,r} \mapsto \mathbb{R}^{p-\ell,r}$ be the identity function and $\bar{\sigma} : \mathbb{R}^{\ell,r} \rightarrow \mathcal{C}(I_r)$ be the function defined for full-rank matrices based on the polar decomposition :

$$\bar{\sigma} : \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} R \in \mathbb{R}^{\ell,r} \mapsto \begin{bmatrix} B_1 \\ B_2 \end{bmatrix},$$

where $B_1 \in \mathcal{S}_r^{++}$, $R \in \mathcal{O}_r$. This definition is correct as the polar decomposition of a full-rank matrix B_1 is unique. Given the orthogonal invariance of F^λ and $F^\lambda \circ \tau$, we have

$$F^\lambda = F_{I_r}^\lambda \circ \tau \circ (\bar{\sigma}, I_{p-\ell,r}) \circ \tau^{-1}.$$

Before applying Theorem 69 with $g = F_{I_r}^\lambda \circ \tau$ and $h = \bar{\sigma} \circ \tau^{-1}$, we first have to prove that its assumptions are satisfied. Clearly, the Jacobian of τ , τ^{-1} and $I_{p-\ell,r}$ are surjective since these are linear invertible functions.

Proposition 70. *Let $A = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \in \mathbb{R}^{\ell,r}$ such that $A_1 \in \mathbb{R}^{r,r}$ is a square invertible matrix and $A_2 \in \mathbb{R}^{\ell-r,r}$. The Jacobian $J\bar{\sigma}(A)$ is a surjective mapping.*

Proof. Thanks to the polar decomposition, we know that there exists $B_1 \in \mathcal{S}_r^{++}$, $B_2 \in \mathbb{R}^{\ell-r,r}$ and $R \in \mathcal{O}_r$ such that

$$A = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix} R.$$

Consequently, we have $\bar{\sigma}(A) = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$.

Also, given $\Delta = \begin{bmatrix} \Delta_1 \\ \Delta_2 \end{bmatrix} \in \mathbb{R}^{\ell,r}$ such that $\Delta_1 \in \mathbb{R}^{r,r}$, $\Delta_2 \in \mathbb{R}^{\ell-r,r}$ and $A + \Delta \in \mathcal{C}_R$, we can write $\begin{bmatrix} \Delta_1 \\ \Delta_2 \end{bmatrix} = \begin{bmatrix} M_1 \\ M_2 \end{bmatrix} R$ where $M_1 \in \mathbb{R}^{r,r}$ is a symmetric matrix, $M_2 \in \mathbb{R}^{\ell-r,r}$ and

$$\bar{\sigma}(A + \Delta) = (A + \Delta)R^T = \bar{\sigma}(A) + \Delta R^T = \bar{\sigma}(A) + \begin{bmatrix} M_1 \\ M_2 \end{bmatrix}.$$

Therefore, we can identify the differential of $\bar{\sigma}$ on the set of matrices $\begin{bmatrix} M_1 \\ M_2 \end{bmatrix} R$ with M_1 symmetric with the linear application $M \mapsto MR^T$. The surjectivity of this differential is obvious. \square

Corollary 71. *Let $0 \leq \lambda < \bar{\lambda}$ and a sublevel set \mathcal{V}^λ be defined as in Corollary 6. The function F^λ has the KL property with exponent $1/2$ in the sublevel set \mathcal{V}^λ .*

Proof. According to Fact 68, $F_{I_r}^\lambda \circ \tau$ is a KL- $\frac{1}{2}$ function around its optimum since it is locally strongly convex. Consequently, $F^\lambda = [F_{I_r}^\lambda \circ \tau] \circ [(\bar{\sigma}, I_{p-\ell,r}) \circ \tau^{-1}]$ is the composition in the sublevel set \mathcal{V}^λ of a KL- $\frac{1}{2}$ function with a smooth function that has a surjective Jacobian mapping, according to Proposition 70. We deduce with Theorem 69 that F^λ has the KL property with exponent $\frac{1}{2}$ in \mathcal{V}^λ . \square

L.2 From KL with exponent $\frac{1}{2}$ to (t-strong proximal PL)

Here, we prove that the KL- $\frac{1}{2}$ property in \mathcal{V}^λ for the function F^λ of SRRR leads to linear convergence in Algorithm 1. This result differs from Theorem 15 of Csiba and Richtárik (2017) for which they assumed strong-convexity instead of the KL property with exponent $\frac{1}{2}$.

As in Theorem 43 we make the assumptions $\mathcal{H}2$ and $\mathcal{H}3$ in this section so that we can use Lemma 46. Indeed, we need these extra assumptions because although the function f we consider for SRRR is L_X -smooth with L_X the largest eigenvalue of $X^T X$, it may not have Lipschitz gradients in the entire sublevel set defined in Corollary 6, mainly because the latter is not convex.

We denote for any $U \in \mathcal{V}^\lambda$, and $t > 0$:

$$\tilde{F}_{t,U}^\lambda(U') := f(U) + \langle \nabla f(U), U' - U \rangle + \frac{1}{2t} \|U' - U\|^2 + \lambda \|U'\|_{1,2}, \quad (110)$$

$$\gamma_t(U) := -\frac{1}{t} \min_{U' \in \mathbb{R}^{p,r}} \left[\tilde{F}_{t,U}^\lambda(U') - F^\lambda(U) \right]. \quad (111)$$

Before obtaining in Proposition 73 a result similar to the (t-strong *proximal* PL) property, we first need to introduce the following result. It is highly similar to Lemma 36 but is adapted to the present context.

Lemma 72. *Let $U \in \mathcal{V}^\lambda$ and $U_+ := \operatorname{argmin}_{U' \in \mathbb{R}^{p,r}} \left[\tilde{F}_{t,U}^\lambda(U') - F^\lambda(U) \right]$. There is a subgradient s_{U_+} of $\|\cdot\|_{1,2}$ at U_+ such that*

$$U_+ - U = -t(\nabla f(U) + \lambda s_{U_+}), \quad (112)$$

$$\gamma_t(U) \geq \frac{1}{2} \|\nabla f(U) + \lambda s_{U_+}\|^2. \quad (113)$$

Proof. Equation (112) is a direct consequence of the first-order optimal conditions for Problem (111). We also have

$$\tilde{F}_{t,U}^\lambda(U_+) - F^\lambda(U) = f(U) + \langle \nabla f(U), U_+ - U \rangle + \frac{1}{2t} \|U_+ - U\|_F^2 + \lambda \|U_+\|_{1,2} - f(U) - \lambda \|U\|_{1,2} \quad (114)$$

$$= \langle \nabla f(U) + \lambda s_{U_+}, U_+ - U \rangle + \frac{t}{2} \|\nabla f(U) + \lambda s_{U_+}\|_F^2 + \lambda [\|U_+\|_{1,2} + \langle s_{U_+}, U - U_+ \rangle - \|U\|_{1,2}], \quad (115)$$

$$\leq -\frac{t}{2} \|\nabla f(U) + \lambda s_{U_+}\|_F^2. \quad (116)$$

In Equation (114), we simply use Equation (110). Equation (115) follows from Equation (112). Equation (116) follows again from Equation (112) and from the convexity of $\|\cdot\|_{1,2}$. Therefore, we have

$$\gamma_t(x) \geq \frac{1}{2} \|\nabla f(x) + \lambda s_{U_+}\|^2.$$

□

Proposition 73. *Let $k_1 \geq 0$ be defined as in Lemma 46, $k \geq k_1$ and assume that $U_k \in \mathcal{V}^\lambda \setminus \Omega^*$. Let $U_{k+1} = \operatorname{argmin}_{U' \in \mathbb{R}^{p,r}} \left[\tilde{F}_{t_k,U_k}^\lambda(U') - F^\lambda(U_k) \right]$. We have*

$$c^2(1 + (Mt_k)^2)\gamma_{t_k}(U_k) \geq F^\lambda(U_{k+1}) - F^{\lambda,*}. \quad (117)$$

Proof. We know from Lemma 72 that there exists a subgradient $s_{U_{k+1}}$ of $U' \mapsto \|U'\|_{1,2}$ at U_{k+1} such that

$$U_{k+1} = U_k - t_k [\nabla f(U_k) + \lambda s_{U_{k+1}}]. \quad (118)$$

We have

$$\|\nabla f(U_{k+1}) + \lambda s_{U_{k+1}}\|^2 \leq 2\|\nabla f(U_k) + \lambda s_{U_{k+1}}\|^2 + 2\|\nabla f(U_k) - \nabla f(U_{k+1})\|^2 \quad (119)$$

$$\leq 2\|\nabla f(U_k) + \lambda s_{U_{k+1}}\|^2 + 2M^2\|U_k - U_{k+1}\|^2 \quad (120)$$

$$\leq 2\|\nabla f(U_k) + \lambda s_{U_{k+1}}\|^2 + 2(Mt_k)^2\|\nabla f(U_k) + \lambda s_{U_{k+1}}\|^2 \quad (121)$$

$$\leq 2(1 + (Mt_k)^2)\|\nabla f(U_k) + \lambda s_{U_{k+1}}\|^2 \quad (122)$$

$$\leq 4(1 + (Mt_k)^2)\gamma_t(U_k) \quad (123)$$

We obtain Equation (119) using the triangle inequality and the inequality of arithmetic and geometric means. Equation (120) is due to Lemma 46. We have Equation (121) thanks to Equation (112). Equation (122) follows from Equation (113) in Lemma 72.

Since $\|\nabla f(U_{k+1}) + \lambda s_{U_{k+1}}\|^2$ is an upper bound of $\text{dist}(0, \partial F^\lambda(U_{k+1}))^2$, Equation 123 implies that

$$\text{dist}(0, \partial F^\lambda(U_{k+1}))^2 \leq 4(1 + (Mt_k)^2)\gamma_{t_k}(U_k).$$

Besides, we know from Corollary 71 that there exists $c > 0$ such that for any $U' \in \mathcal{V}^\lambda$, the function F^λ satisfies the inequality :

$$\text{dist}(0, \partial F^\lambda(U')) \geq \frac{2}{c} \sqrt{F^\lambda(U') - F^{\lambda,*}}. \quad (124)$$

The element U_{k+1} being in the sublevel set \mathcal{V}^λ since $F^\lambda(U_{k+1}) \leq F^\lambda(U_k)$ and $U_k \in \mathcal{V}^\lambda$, we finally obtain with Equation (124) :

$$(1 + (Mt_k)^2)\gamma_{t_k}(U_k) \geq \frac{1}{c^2}(F^\lambda(U_{k+1}) - F^{\lambda,*})$$

□

Remark 74. Note that Equation (t-strong proximal PL) in Section 5.3, that is to say in the PL framework, can be written

$$\gamma_t(U) \geq c_1[F^\lambda(U) - F^{\lambda,*}] \quad \text{with } c_1 > 0,$$

while Equation 117, in the KL framework, can be written,

$$\gamma_t(U) \geq c_2[F^\lambda(U_+) - F^{\lambda,*}] \quad \text{with } c_2 > 0.$$

The right term depends either on U or U_+ and this is the main reason for the differences found in the computations between the two frameworks.

We denote

$$\xi : U' \mapsto F^\lambda(U') - F^{\lambda,*}.$$

Proposition 73 finally leads to local linear convergence, as encountered in Proposition 5.1 of Li and Pong (2017) for batch proximal gradient descent. As in Proposition 5.1 of Li and Pong (2017), we have to use an upper bound $d > 0$ on the step size t while this was not necessary when we used the Polyak-Łojasiewicz inequality instead of the Kurdyka-Łojasiewicz inequality.

Proposition 75. Let $k_1 \geq 0$ be defined as in Lemma 46. Assume that there is $d > 0$ such that for any $k \geq k_1$, we have $t_k \leq d$. There is $0 < \rho < 1$ such that for any $k \geq k_1$, if $U_k \in \mathcal{V}^\lambda \setminus \Omega^*$, then we have :

$$\xi(U_{k+1}) \leq (1 - \rho)\xi(U_k).$$

Therefore, the convergence of Algorithm 1 is locally linear.

Proof. Let $k \geq k_1$, $U_k \in \mathcal{V}^\lambda$ and

$$U_{k+1} = \underset{U' \in \mathbb{R}^{p,r}}{\text{argmin}} \left[\tilde{F}_{t_k, U_k}^\lambda(U') - F^\lambda(U_k) \right].$$

First, from Equation (117) in Proposition 73, we have

$$\frac{\gamma_{t_k}(U_k)}{\xi(U_{k+1})} \geq \frac{1}{c^2(1 + (Mt_k)^2)}. \quad (125)$$

Secondly, we have

$$\xi(U_{k+1}) \leq \xi(U_k) - t_k \gamma_{t_k}(U_k) \quad (126)$$

$$\begin{aligned} &\leq \xi(U_k) - t_k \frac{\gamma_t(U_k)}{\xi(U_{k+1})} \xi(U_{k+1}) \\ &\leq \xi(U_k) - t_k \frac{1}{c^2(1+(Mt_k)^2)} \xi(U_{k+1}) \end{aligned} \quad (127)$$

$$\leq \xi(U_k) - t_k \frac{1}{c^2(1+(Md)^2)} \xi(U_{k+1}). \quad (128)$$

Equation (126) comes from Fact 33. Equation (127) follows from Equation 125 and Equation (128) follows from the assumption $t_k \leq d$. Consequently, we have

$$\begin{aligned} \xi(U_{k+1}) &\leq \frac{1}{1 + \frac{t_k}{c^2(1+(Md)^2)}} \xi(U_k) \\ &\leq \frac{1}{1 + \frac{\beta}{c^2 L_X(1+(Md)^2)}} \xi(U_k) \\ &\leq (1 - \rho) \xi(U_k) \quad \text{with } \rho := 1 - \frac{1}{1 + \frac{\beta}{c^2 L_X(1+(Md)^2)}}. \end{aligned} \quad (129)$$

We have Inequality (129) since $t_k > \frac{\beta}{L_X}$ for k sufficiently large by Fact 35. \square

Proposition 75 finally leads to local linear convergence for the forward-backward algorithm applied to (SRRR). The proof in this section is different from the core of the article since we used KL inequalities instead of PL inequalities for the proof.

M Additional details and results on the experiments

M.1 Algorithm of Park et al. (2016)

To evaluate the performance of Algorithm 1 for RRR, we compare it with the algorithm proposed in Park et al. (2016), which minimizes the biconvex formulation of Problem (1) *i.e.* with the form of Equation (4). To avoid the scaling issue due to the invariance of the objective by any transformation $(U, V) \mapsto (UC, VC^{-T})$ where C is a square invertible matrix, the formulation they propose to add a regularizer $(U, V) \mapsto \frac{1}{4} \|U^T U - V^T V\|_F^2$ which does not change the optimal value of the function. With this differentiable function, simultaneous gradient descent in U and V is feasible. However, this regularization scheme is not applicable if a group-Lasso penalty is added, because the latter is not compatible with imposing the constraint $U^T U = V^T V$ at the optimum.

M.2 Different values of the correlation coefficient ρ

Given that the choice of ρ has a strong impact on the running time, we report in Figure 5 and Figure 6 additional results for different values of the parameter ρ . Apart from this modification, we test the algorithms with the same setting as in Section 6. This change corresponds to modifying the correlation between the columns of the design matrix X . Although the speed of the algorithms decreases when ρ increases, the relative order of the methods remains the same.

M.3 Different sparsity scenarios

To assess the quality of the algorithm when the proportion of zero rows in W_0 varies, we present Figure 7 where the proportion of zero rows p_0 is respectively 0.5 and 0.8, that is W_0 has 50% and 80% of zero rows.

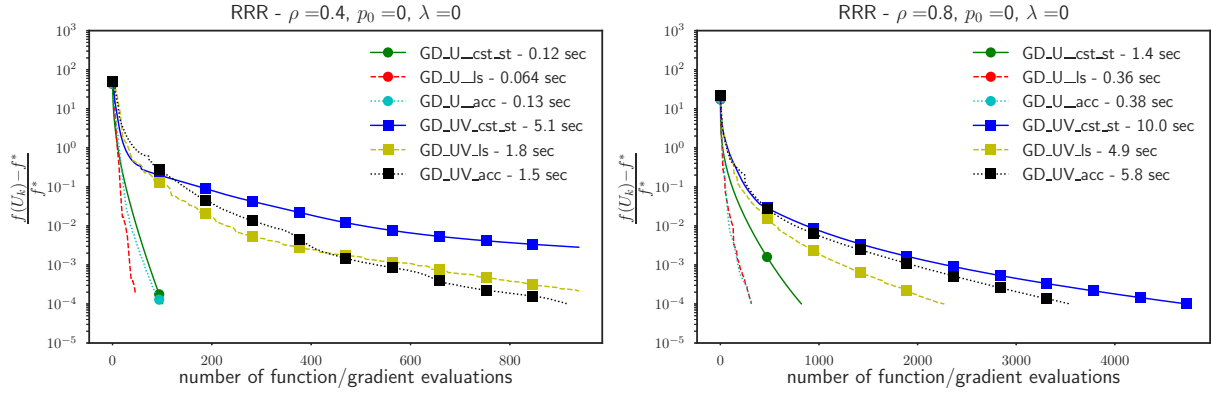


Figure 5: (Left) RRR : $\rho = 0.4$. (Right) RRR : $\rho = 0.8$. Times reported are times to reach a gap of 10^{-4}

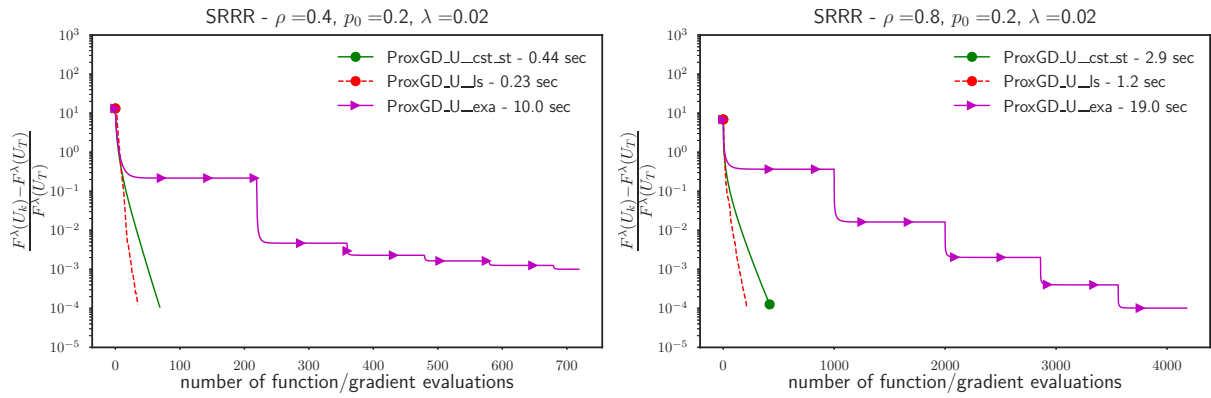


Figure 6: (Left) SRRR : $\rho = 0.4$. (Right) SRRR : $\rho = 0.8$. Times reported are times to reach a gap of 10^{-4} .

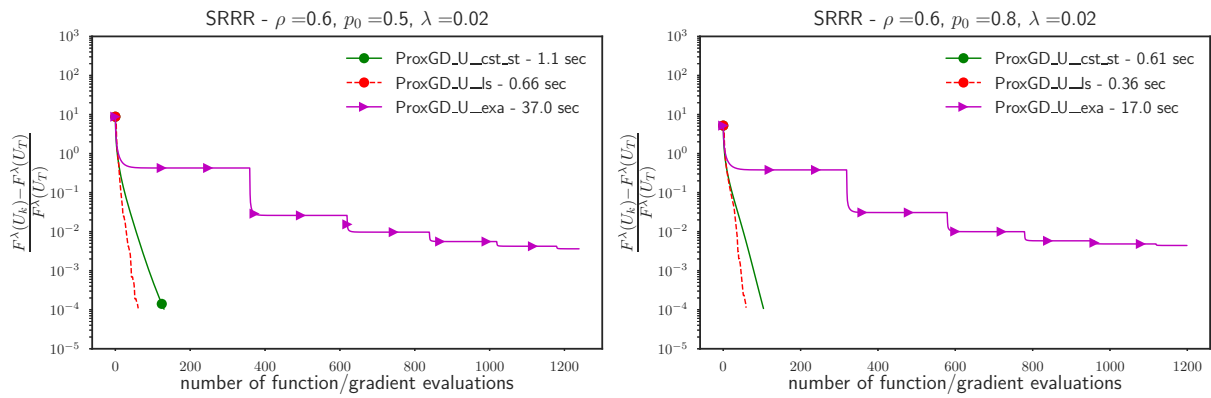


Figure 7: (Left) SRRR : $\rho = 0.6, p_0 = 0.5$ and $\lambda = 0.02$. (Right) SRRR : $\rho = 0.6, p_0 = 0.8$ and $\lambda = 0.02$. Times reported are times to reach a gap of 10^{-4} .

References

- Attouch, H. and Bolte, J. (2009). On the convergence of the proximal algorithm for nonsmooth functions involving analytic features. *Mathematical Programming*, 116(1-2):5–16.
- Attouch, H., Bolte, J., Redont, P., and Soubeyran, A. (2010). Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-łojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457.
- Attouch, H., Bolte, J., and Svaiter, B. F. (2013). Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized gauss-seidel methods. *Mathematical Programming*, 137(1-2):91–129.
- Bach, F., Jenatton, R., Mairal, J., Obozinski, G., et al. (2012). Optimization with sparsity-inducing penalties. *Foundations and Trends® in Machine Learning*, 4(1):1–106.
- Bolte, J., Daniilidis, A., and Lewis, A. (2007). The łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. *SIAM Journal on Optimization*, 17(4):1205–1223.
- Bonnans, J. F. and Shapiro, A. (1998). Optimization problems with perturbations: A guided tour. *SIAM review*, 40(2):228–264.
- Chouzenoux, E., Pesquet, J.-C., and Repetti, A. (2014). Variable metric forward-backward algorithm for minimizing the sum of a differentiable function and a convex function. *Journal of Optimization Theory and Applications*, 162(1):107–132.
- Csiba, D. and Richtárik, P. (2017). Global convergence of arbitrary-block gradient methods for generalized Polyak-Łojasiewicz functions. *arXiv preprint arXiv:1709.03014*.
- Danskin, J. M. (1967). *The theory of max-min and its application to weapons allocation problems*, volume 5. Springer Science & Business Media.
- Frankel, P., Garrigos, G., and Peypouquet, J. (2015). Splitting methods with variable metric for Kurdyka-Łojasiewicz functions and general convergence rates. *Journal of Optimization Theory and Applications*, 165(3):874–900.
- Ge, R., Jin, C., and Zheng, Y. (2017). No spurious local minima in nonconvex low rank problems: a unified geometric analysis. *arXiv preprint arXiv:1704.00708*.
- Grave, E., Obozinski, G. R., and Bach, F. R. (2011). Trace lasso: a trace norm regularization for correlated designs. In *Advances in Neural Information Processing Systems*, pages 2187–2195.
- Khamaru, K. and Wainwright, M. (2018). Convergence guarantees for a class of non-convex and non-smooth optimization problems. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2601–2610.
- Li, G. and Pong, T. K. (2017). Calculus of the exponent of Kurdyka-Łojasiewicz inequality and its applications to linear convergence of first-order methods. *Foundations of Computational Mathematics*, pages 1–34.
- Li, X., Wang, Z., Lu, J., Arora, R., Haupt, J., Liu, H., and Zhao, T. (2016). Symmetry, saddle points, and global geometry of nonconvex matrix factorization. *arXiv preprint arXiv:1612.09296*.
- Oliver, H. W. (1954). The exact Peano derivative. *Transactions of the American Mathematical Society*, 76(3):444–456.
- Park, D., Kyrillidis, A., Caramanis, C., and Sanghavi, S. (2016). Finding low-rank solutions via non-convex matrix factorization, efficiently and provably. *arXiv preprint arXiv:1606.03168*.

Rockafellar, R. T. and Wets, R. J.-B. (2009). *Variational analysis*, volume 317. Springer Science & Business Media.

Stewart, G. (2012). Smooth local bases for perturbed eigenspaces. *Institute for Advanced Computer Studies TR*, page 08.