

## A. Auxiliary Lemmas

In this section, we prove Lemma A.1 and a few auxiliary lemmas that we will need for the proofs of Theorem 2.4 and Theorem 3.6.

**Lemma A.1.** Let  $\mathbf{x} \in \mathbb{R}^{d_2}$  be distributed according to distribution  $\mathcal{D}$  with  $\mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}$ . Then, for  $\ell(\mathbf{U}, \mathbf{V}) := \mathbb{E}_{\mathbf{x}}[\|\mathbf{y} - \mathbf{UV}^\top \mathbf{x}\|^2]$  and  $f(\mathbf{U}, \mathbf{V}) := \mathbb{E}_{\mathbf{b}, \mathbf{x}}[\|\mathbf{y} - \frac{1}{\theta} \mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V}^\top \mathbf{x}\|^2]$ , it holds that

$$f(\mathbf{U}, \mathbf{V}) = \ell(\mathbf{U}, \mathbf{V}) + \lambda \sum_{i=1}^r \|\mathbf{u}_i\|^2 \|\mathbf{v}_i\|^2. \quad (8)$$

Furthermore,  $\ell(\mathbf{U}, \mathbf{V}) = \|\mathbf{M} - \mathbf{UV}^\top\|_F^2$ .

*Proof of Lemma A.1.* The proof closely follows (Cavazza et al., 2018). Recall that  $\mathbf{y} = \mathbf{M}\mathbf{x}$ , for some unknown  $\mathbf{M} \in \mathbb{R}^{d_2 \times d_1}$ . Observe that

$$\begin{aligned} f(\mathbf{U}, \mathbf{V}) &= \mathbb{E}_{\mathbf{x}}[\|\mathbf{y}\|^2] + \frac{1}{\theta^2} \mathbb{E}_{\mathbf{b}, \mathbf{x}}[\|\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V}^\top \mathbf{x}\|^2] \\ &\quad - \frac{2}{\theta} \mathbb{E}_{\mathbf{x}}[\langle \mathbf{M}\mathbf{x}, \mathbb{E}_{\mathbf{b}}[\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V}^\top \mathbf{x}] \rangle] \end{aligned} \quad (9)$$

where we used the fact that  $\mathbf{y} = \mathbf{M}\mathbf{x}$ . We have the following set of equalities for the second term on the right hand side of Equation (9):

$$\begin{aligned} \mathbb{E}_{\mathbf{b}, \mathbf{x}}[\|\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V}^\top \mathbf{x}\|^2] &= \mathbb{E}_{\mathbf{x}} \sum_{i=1}^{d_2} \mathbb{E}_{\mathbf{b}} \left( \sum_{j=1}^r u_{ij} b_j v_j^\top \mathbf{x} \right)^2 \\ &= \mathbb{E}_{\mathbf{x}} \sum_{i=1}^{d_2} \mathbb{E}_{\mathbf{b}} \left[ \sum_{j,k=1}^r u_{ij} u_{ik} b_j b_k (v_j^\top \mathbf{x})(v_k^\top \mathbf{x}) \right] \\ &= \mathbb{E}_{\mathbf{x}} \sum_{i=1}^{d_2} \sum_{j,k=1}^r u_{ij} u_{ik} (\theta^2 1_{j \neq k} + \theta 1_{j=k}) (v_j^\top \mathbf{x})(v_k^\top \mathbf{x}) \\ &= \theta^2 \mathbb{E}_{\mathbf{x}}[\|\mathbf{UV}^\top \mathbf{x}\|^2] + (\theta - \theta^2) \mathbb{E}_{\mathbf{x}} \sum_{i=1}^{d_2} \sum_{j=1}^r u_{ij}^2 (v_j^\top \mathbf{x})^2 \\ &= \theta^2 \mathbb{E}_{\mathbf{x}}[\|\mathbf{UV}^\top \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^r \|\mathbf{v}_j\|^2 \sum_{i=1}^{d_2} u_{ij}^2 \\ &= \theta^2 \mathbb{E}_{\mathbf{x}}[\|\mathbf{UV}^\top \mathbf{x}\|^2] + (\theta - \theta^2) \sum_{j=1}^r \|\mathbf{v}_j\|^2 \|\mathbf{u}_j\|^2, \end{aligned} \quad (10)$$

where the second to last equality follows because  $\mathbb{E}_{\mathbf{x}}[(v_j^\top \mathbf{x})^2] = v_j^\top \mathbb{E}_{\mathbf{x}}[\mathbf{x}\mathbf{x}^\top] v_j = \|\mathbf{v}_j\|^2$ . For the third term in Equation (9) we have:

$$\langle \mathbf{M}\mathbf{x}, \mathbb{E}_{\mathbf{b}}[\mathbf{U} \text{diag}(\mathbf{b}) \mathbf{V}^\top \mathbf{x}] \rangle = \theta \langle \mathbf{M}\mathbf{x}, \mathbf{UV}^\top \mathbf{x} \rangle \quad (11)$$

Plugging Equations (10) and (11) into (9), we get

$$\begin{aligned} f(\mathbf{U}, \mathbf{V}) &= \mathbb{E}_{\mathbf{x}}[\|\mathbf{y}\|^2] + \mathbb{E}_{\mathbf{x}}[\|\mathbf{UV}^\top \mathbf{x}\|^2] - 2\mathbb{E}_{\mathbf{x}} \langle \mathbf{M}\mathbf{x}, \mathbf{UV}^\top \mathbf{x} \rangle \\ &\quad + \frac{1-\theta}{\theta} \sum_{i=1}^r \|\mathbf{u}_i\|^2 \|\mathbf{v}_i\|^2 \end{aligned} \quad (12)$$

It is easy to check that the first three terms in Equation (12) sum to  $\ell(\mathbf{U}, \mathbf{V})$ . Furthermore, since for any  $\mathbf{A} \in \mathbb{R}^{d_2 \times d_1}$  it holds that  $\|\mathbf{A}\mathbf{x}\|^2 = \|\mathbf{A}\|_F^2$ , we should have  $\ell(\mathbf{U}, \mathbf{V}) = \|\mathbf{M} - \mathbf{UV}^\top\|_F^2$ .  $\square$

**Lemma A.2.** For any pair of integers  $\rho$  and  $r$ , and for any  $\lambda \in \mathbb{R}_+$ , it holds that

$$(\mathbf{I}_\rho + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^\top)^{-1} = \mathbf{I}_\rho - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^\top.$$

Lemma A.2 is an instance of the Woodbury's matrix identity. Here, we include a proof for completeness.

*Proof of Lemma A.2.* The proof simply follows from the following set of equations.

$$\begin{aligned} &(\mathbf{I}_\rho + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^\top)(\mathbf{I}_\rho - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^\top) \\ &= \mathbf{I}_\rho + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^\top - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^\top - \frac{\lambda^2}{r(r + \lambda\rho)} \mathbf{1}\mathbf{1}^\top \mathbf{1}\mathbf{1}^\top \\ &= \mathbf{I}_\rho + \left( \frac{\lambda}{r} - \frac{\lambda}{r + \lambda\rho} - \frac{\rho\lambda^2}{r(r + \lambda\rho)} \right) \mathbf{1}\mathbf{1}^\top = \mathbf{I}_\rho \end{aligned}$$

$\square$

**Lemma A.3.** Let  $\lambda > 0$  be a constant. Let  $\mathbf{a} \in \mathbb{R}_+^d$  such that  $a_i \geq a_{i+1}$  for all  $i \in [d-1]$ . For  $r \leq d$ , let the function  $g: [r] \rightarrow \mathbb{R}$  be defined as

$$\begin{aligned} g(\rho) &:= \sum_{i=1}^{\rho} \left( \frac{\lambda \sum_{k=1}^{\rho} a_k}{r + \lambda\rho} \right)^2 + \sum_{i=\rho+1}^d a_i^2 \\ &\quad + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \left( a_i - \frac{\lambda \sum_{k=1}^{\rho} a_k}{r + \lambda\rho} \right) \right)^2. \end{aligned}$$

Then  $g(\rho)$  is monotonically non-increasing in  $\rho$ .

*Proof of Lemma A.3.* Let denote the sum of the top  $\tau$  elements of  $\mathbf{a}$  by  $h_\tau = \sum_{i=1}^{\tau} a_i$ . Furthermore, let the sum of squared of  $\tau$  bottom elements of  $\mathbf{a}$  be denoted by  $t_\tau = \sum_{i=\tau+1}^d a_i^2$ . We can simplify  $g(\rho)$  and give it in terms of  $h_\rho$  and  $t_\rho$  as follows:

$$\begin{aligned} g(\rho) &= \rho \left( \frac{\lambda h_\rho}{r + \lambda\rho} \right)^2 + t_\rho + \frac{\lambda}{r} \left( \left( 1 - \frac{\lambda\rho}{r + \lambda\rho} \right) h_\rho \right)^2 \\ &= \frac{\rho\lambda^2 + \lambda r}{(r + \lambda\rho)^2} (h_\rho)^2 + t_\rho \\ &= \frac{\lambda h_\rho^2}{r + \lambda\rho} + t_\rho \end{aligned}$$

It suffices to show that  $g(\rho + 1) \leq g(\rho)$  for all  $\rho \in [r - 1]$ .  $\square$

$$\begin{aligned}
 g(\rho + 1) &= \frac{\lambda h_{\rho+1}^2}{r + \lambda\rho + \lambda} + t_{\rho+1} \\
 &= \frac{\lambda}{r + \lambda\rho + \lambda} (h_\rho^2 + \lambda_{\rho+1}^2(\mathbf{M}) + 2\lambda_{\rho+1}(\mathbf{M})h_\rho) \\
 &\quad - \lambda_{\rho+1}^2(\mathbf{M}) + t_\rho \\
 &= g(\rho) - \frac{\lambda^2 h_\rho^2}{(r + \lambda\rho)(r + \lambda\rho + \lambda)} - \lambda_{\rho+1}^2(\mathbf{M}) \\
 &\quad + \frac{\lambda}{r + \lambda\rho + \lambda} (\lambda_{\rho+1}^2(\mathbf{M}) + 2\lambda_{\rho+1}(\mathbf{M})h_\rho) \\
 &= g(\rho) - \frac{\lambda^2 h_\rho^2}{(r + \lambda\rho)(r + \lambda\rho + \lambda)} - \frac{(r + \lambda\rho)\lambda_{\rho+1}^2(\mathbf{M})}{r + \lambda\rho + \lambda} \\
 &\quad + \frac{\lambda}{r + \lambda\rho + \lambda} (2\lambda_{\rho+1}(\mathbf{M})h_\rho) \\
 &= g(\rho) - \frac{(\lambda h_\rho - (r + \lambda\rho)\lambda_{\rho+1}^2(\mathbf{M}))^2}{(r + \lambda\rho)(r + \lambda\rho + \lambda)} \leq g(\rho).
 \end{aligned}$$

Hence  $g(\rho)$  is monotonically non-increasing in  $\rho$ .  $\square$

## B. Proofs of Theorems in Section 2

*Proof of Theorem 2.2.* Consider the matrix  $\mathbf{G}_1 := \mathbf{G}_U - \frac{\text{Tr} \mathbf{G}_U}{r} \mathbf{I}_r$ . We exhibit an orthogonal transformation  $\mathbf{Q}$ , such that  $\mathbf{Q}^\top \mathbf{G}_1 \mathbf{Q}$  is zero on its diagonal. Observe that

$$\mathbf{Q}^\top \mathbf{G}_U \mathbf{Q} = \mathbf{Q}^\top \mathbf{G}_1 \mathbf{Q} + \frac{\text{Tr} \mathbf{G}_U}{r} \mathbf{I}_r,$$

so that all diagonal elements of  $\mathbf{G}_U$  are equal to  $\frac{\text{Tr} \mathbf{G}_U}{r}$ , i.e.  $\mathbf{G}_U$  is equalized.

Our construction closely follows the proof of a classical theorem in matrix analysis, which states that any trace zero matrix is a commutator (Albert and Muckenhoupt, 1957; Kahan, 1999). For the zero trace matrix  $\mathbf{G}_1$ , we first show that there exists a unit vector  $\mathbf{w}_{11}$  such that  $\mathbf{w}_{11}^\top \mathbf{G}_1 \mathbf{w}_{11} = 0$ .

**Claim 1.** Assume  $\mathbf{G}$  is a zero trace matrix and let  $\mathbf{G} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$  be an eigendecomposition of  $\mathbf{G}$ . Then  $\mathbf{w} = \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{u}_i$  has a vanishing Rayleigh quotient, that is,  $\mathbf{w}^\top \mathbf{G} \mathbf{w} = 0$ , and  $\|\mathbf{w}\| = 1$ .

*Proof of Claim 1.* First, we notice that  $\mathbf{w}$  has unit norm

$$\|\mathbf{w}\|^2 = \left\| \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{u}_i \right\|^2 = \frac{1}{r} \left\| \sum_{i=1}^r \mathbf{u}_i \right\|^2 = \frac{1}{r} \sum_{i=1}^r \|\mathbf{u}_i\|^2 = 1.$$

It is easy to see that  $\mathbf{w}$  has a zero Rayleigh quotient

$$\begin{aligned}
 \mathbf{w}^\top \mathbf{G} \mathbf{w} &= \left( \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{u}_i \right)^\top \mathbf{G} \left( \frac{1}{\sqrt{r}} \sum_{i=1}^r \mathbf{u}_i \right) \\
 &= \frac{1}{r} \sum_{i,j=1}^r \mathbf{u}_i^\top \mathbf{G} \mathbf{u}_j = \frac{1}{r} \sum_{i=1}^r \lambda_j \mathbf{u}_i^\top \mathbf{u}_j = \frac{1}{r} \sum_{i=1}^r \lambda_i = 0.
 \end{aligned}$$

Let  $\mathbf{W}_1 := [\mathbf{w}_{11}, \mathbf{w}_{12}, \dots, \mathbf{w}_{1d}]$  be such that  $\mathbf{W}_1^\top \mathbf{W}_1 = \mathbf{W}_1 \mathbf{W}_1^\top = \mathbf{I}_d$ . Observe that  $\mathbf{W}_1^\top \mathbf{G}_1 \mathbf{W}_1$  has zero on its first diagonal elements

$$\mathbf{W}_1^\top \mathbf{G}_1 \mathbf{W}_1 = \begin{bmatrix} 0 & \mathbf{b}_1^\top \\ \mathbf{b}_1 & \mathbf{G}_2 \end{bmatrix}$$

The principal submatrix  $\mathbf{G}_2$  also has a zero trace. With a similar argument, let  $\mathbf{w}_{22} \in \mathbb{R}^{d-1}$  be such that  $\|\mathbf{w}_{22}\| = 1$  and  $\mathbf{w}_{22}^\top \mathbf{G}_2 \mathbf{w}_{22} = 0$  and define

$\mathbf{W}_2 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \mathbf{w}_{22} & \mathbf{w}_{23} & \dots & \mathbf{w}_{2d} \end{bmatrix} \in \mathbb{R}^{d \times d}$  such that  $\mathbf{W}_2^\top \mathbf{W}_2 = \mathbf{W}_2 \mathbf{W}_2^\top = \mathbf{I}_d$ , and observe that

$$(\mathbf{W}_1 \mathbf{W}_2)^\top \mathbf{G}_1 (\mathbf{W}_1 \mathbf{W}_2) = \begin{bmatrix} 0 & \cdot & \dots \\ \cdot & 0 & \dots \\ \vdots & \vdots & \mathbf{G}_2 \end{bmatrix}.$$

This procedure can be applied recursively so that for the equalizer  $\mathbf{Q} = \mathbf{W}_1 \mathbf{W}_2 \dots \mathbf{W}_d$  we have

$$\mathbf{Q}^\top \mathbf{G}_1 \mathbf{Q} = \begin{bmatrix} 0 & \cdot & \dots & \cdot \\ \cdot & 0 & \dots & \cdot \\ \vdots & \vdots & \ddots & \vdots \\ \cdot & \cdot & \cdot & 0 \end{bmatrix}.$$

$\square$

*Proof of Theorem 2.3.* Let us denote the squared column norms of  $\mathbf{U}$  by  $\mathbf{n}_u = (\|\mathbf{u}_1\|^2, \dots, \|\mathbf{u}_r\|^2)$ . Observe that for any weight matrix  $\mathbf{U}$ :

$$\begin{aligned}
 R(\mathbf{U}, \mathbf{U}) &= \lambda \sum_{i=1}^r \|\mathbf{u}_i\|^4 = \frac{\lambda}{r} \|\mathbf{1}_r\|^2 \|\mathbf{n}_u\|^2 \\
 &\geq \frac{\lambda}{r} \langle \mathbf{1}_r, \mathbf{n}_u \rangle^2 = \frac{\lambda}{r} \left( \sum_{i=1}^r \|\mathbf{u}_i\|^2 \right)^2 = \frac{\lambda}{r} \|\mathbf{U}\|_F^4,
 \end{aligned}$$

where  $\mathbf{1}_r \in \mathbb{R}^r$  is the vector of all ones and the inequality is due to Cauchy-Schwartz. Hence, the regularizer is lower bounded by  $\frac{\lambda}{r} \|\mathbf{U}\|_F^4$ , with equality if and only if  $\mathbf{n}_u$  is parallel to  $\mathbf{1}_r$ , i.e. when  $\mathbf{U}$  is equalized. Now, if  $\mathbf{U}$  is not equalized, by Theorem 2.2 there exist a rotation matrix  $\mathbf{Q}$  such that  $\mathbf{U}\mathbf{Q}$  is equalized, which implies  $R(\mathbf{U}\mathbf{Q}, \mathbf{U}\mathbf{Q}) < R(\mathbf{U}, \mathbf{U})$ . Together with rotational invariance of the loss function, this gives a contradiction with global optimality  $\mathbf{U}$ . Hence, if  $\mathbf{U}$  is a global optimum then it is equalized and we have  $R(\mathbf{U}, \mathbf{U}) = \lambda \sum_{i=1}^r \|\mathbf{u}_i\|^4 = \frac{\lambda}{r} \|\mathbf{U}\|_F^4$ .  $\square$

*Proof of Theorem 2.4.* By Theorem 2.3, if  $\mathbf{W}$  is an optimum of Problem 4, then it holds that  $\lambda \sum_{i=1}^r \|\mathbf{w}_i\|^4 = \frac{\lambda}{r} \|\mathbf{W}\|_F^4$ . Also, by Theorem 2.2, it is always possible to equalize

any given weight matrix. Hence, Problem 4 reduces to the following problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times r}} \|\mathbf{M} - \mathbf{W}\mathbf{W}^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{W}\|_F^4 \quad (13)$$

Let  $\mathbf{M} = \mathbf{U}_M \Lambda_M \mathbf{U}_M^\top$  and  $\mathbf{W} = \mathbf{U}_W \Sigma_W \mathbf{V}_W^\top$  be an eigen-decomposition of  $\mathbf{M}$  and a full SVD of  $\mathbf{W}$  respectively, such that  $\lambda_i(\mathbf{M}) \geq \lambda_{i+1}(\mathbf{M})$  and  $\sigma_i(\mathbf{W}) \geq \sigma_{i+1}(\mathbf{W})$  for all  $i \in [d-1]$ . Rewriting objective of Problem 13 in terms of these decompositions gives:

$$\begin{aligned} & \|\mathbf{M} - \mathbf{W}\mathbf{W}^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{W}\|_F^4 \\ &= \|\mathbf{U}_M \Lambda_M \mathbf{U}_M^\top - \mathbf{U}_W \Sigma_W \Sigma_W^\top \mathbf{U}_W^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{U}_W \Sigma_W \mathbf{V}_W^\top\|_F^4 \\ &= \|\Lambda_M - \mathbf{U}' \Sigma_W \Sigma_W^\top \mathbf{U}'^\top\|_F^2 + \frac{\lambda}{r} \|\Sigma_W\|_F^4 \\ &= \|\Lambda_M\|_F^2 + \|\Lambda_W\|_F^2 - 2\langle \Lambda_M, \mathbf{U}' \Lambda_W \mathbf{U}'^\top \rangle + \frac{\lambda}{r} (\text{Tr}(\Lambda_W))^2 \end{aligned}$$

where  $\Lambda_W := \Sigma_W \Sigma_W^\top$  and  $\mathbf{U}' = \mathbf{U}_M^\top \mathbf{U}_W$ . By Von Neumann's trace inequality, for a fixed  $\Sigma_W$  we have that

$$\langle \Lambda_M, \mathbf{U}' \Lambda_W \mathbf{U}'^\top \rangle \leq \sum_{i=1}^d \lambda_i(\mathbf{M}) \lambda_i(\mathbf{W}),$$

where the equality is achieved when  $\Lambda_i(\mathbf{W})$  have the same ordering as  $\Lambda_i(\mathbf{M})$  and  $\mathbf{U}' = \mathbf{I}$ , i.e.  $\mathbf{U}_M = \mathbf{U}_W$ . Now, Problem 13 is reduced to

$$\begin{aligned} & \min_{\substack{\|\Lambda_W\|_0 \leq r, \\ \Lambda_W \geq 0}} \|\Lambda_M - \Lambda_W\|_F^2 + \frac{\lambda}{r} (\text{Tr}(\Lambda_W))^2 \\ &= \min_{\bar{\lambda} \in \mathbb{R}_+^r} \sum_{i=1}^r (\lambda_i(\mathbf{M}) - \bar{\lambda}_i)^2 + \sum_{i=r+1}^d \lambda_i^2(\mathbf{M}) + \frac{\lambda}{r} \left( \sum_{i=1}^r \bar{\lambda}_i \right)^2 \end{aligned}$$

The Lagrangian is given by

$$\begin{aligned} L(\bar{\lambda}, \alpha) &= \sum_{i=1}^r (\lambda_i(\mathbf{M}) - \bar{\lambda}_i)^2 + \sum_{i=r+1}^d \lambda_i^2(\mathbf{M}) \\ &\quad + \frac{\lambda}{r} \left( \sum_{i=1}^r \bar{\lambda}_i \right)^2 - \sum_{i=1}^r \alpha_i \bar{\lambda}_i \end{aligned}$$

The KKT conditions ensures that at the optima it holds for all  $i \in [r]$  that

$$\begin{aligned} & \bar{\lambda}_i \geq 0, \alpha_i \geq 0, \bar{\lambda}_i \alpha_i = 0 \\ & 2(\bar{\lambda}_i - \lambda_i(\mathbf{M})) + \frac{2\lambda}{r} \left( \sum_{i=1}^r \bar{\lambda}_i \right) - \alpha_i = 0 \end{aligned}$$

Let  $\rho = |\{i : \bar{\lambda}_i > 0\}| \leq r$  be the number of nonzero  $\bar{\lambda}_i$ . For

$i = 1, \dots, \rho$  we have  $\alpha_i = 0$ , hence

$$\begin{aligned} & \bar{\lambda}_i + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \bar{\lambda}_i \right) = \lambda_i(\mathbf{M}) \\ & \implies (\mathbf{I}_\rho + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^\top) \bar{\lambda}_{1:\rho} = \lambda_{1:\rho}(\mathbf{M}) \\ & \implies \bar{\lambda}_{1:\rho} = (\mathbf{I}_\rho - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^\top) \lambda_{1:\rho}(\mathbf{M}) \\ & \implies \bar{\lambda}_{1:\rho} = \lambda_{1:\rho}(\mathbf{M}) - \frac{\lambda\rho\kappa_\rho}{r + \lambda\rho} \mathbf{1}_\rho \\ & \implies \Lambda_W = (\Lambda_M - \frac{\lambda\rho\kappa_\rho}{r + \lambda\rho} \mathbf{I}_d)_+ \end{aligned}$$

where  $\kappa_\rho := \frac{1}{\rho} \sum_{i=1}^{\rho} \lambda_i(\mathbf{M})$  and the second implication is due to Lemma A.2. It only remains to find the optimal  $\rho$ . Let's define the function

$$\begin{aligned} g(\rho) &:= \sum_{i=1}^{\rho} (\lambda_i(\mathbf{M}) - \bar{\lambda}_i)^2 + \sum_{i=\rho+1}^d \lambda_i^2(\mathbf{M}) + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \bar{\lambda}_i \right)^2 \\ &= \sum_{i=1}^{\rho} \left( \frac{\lambda \sum_{k=1}^{\rho} \lambda_k(\mathbf{M})}{r + \lambda\rho} \right)^2 + \sum_{i=\rho+1}^d \lambda_i(\mathbf{M})^2 \\ &\quad + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \left( \lambda_i(\mathbf{M}) - \frac{\lambda \sum_{k=1}^{\rho} \lambda_k(\mathbf{M})}{r + \lambda\rho} \right) \right)^2. \end{aligned}$$

By Lemma A.3,  $g(\rho)$  is monotonically non-increasing in  $\rho$ , hence  $\rho$  should be the largest *feasible* integer, i.e.

$$\rho = \max\{j : \lambda_j > \frac{\lambda_j \kappa}{r + \lambda_j}\}.$$

□

*Proof of Remark 5.2.* For  $\tilde{\mathbf{U}}$  to have equal column norms, it suffices to show that  $\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}}$  is constant on its diagonal. Next, we note that

$$\begin{aligned} \tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} &= \mathbf{Q}^\top \mathbf{U}^\top \mathbf{U} \mathbf{Q} \\ &= (\mathbf{V} \mathbf{Z}_k)^\top (\mathbf{W} \Sigma \mathbf{V}^\top)^\top (\mathbf{W} \Sigma \mathbf{V}^\top) (\mathbf{V} \mathbf{Z}_k) \\ &= \mathbf{Z}_k^\top \mathbf{V}^\top \mathbf{V} \Sigma \mathbf{W}^\top \mathbf{W} \Sigma \mathbf{V}^\top \mathbf{V} \mathbf{Z}_k \\ &= \mathbf{Z}_k^\top \Sigma^2 \mathbf{Z}_k \end{aligned}$$

It remains to show that for any diagonal matrix  $\mathbf{D}$ ,  $\mathbf{Z}_k^\top \mathbf{D} \mathbf{Z}_k$  is diagonalized. First note that

$$\mathbf{Z}_2 \mathbf{Z}_2^\top = \frac{1}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} = \mathbf{I}_2$$

so that  $\mathbf{Z}_2$  is indeed a rotation. By induction, it is easy to see that  $\mathbf{Z}_k$  is a rotation for all  $k$ . Now, we show that  $\mathbf{Z}_k$  equalizes any diagonal matrix  $\mathbf{D}$ . Observe that

$$[\mathbf{Z}_k^\top \mathbf{D} \mathbf{Z}_k]_{ii} = \sum_{j=1}^{2^{k-1}} D_{jj} z_{ji}^2 = \sum_{j=1}^{2^{k-1}} D_{ii} 2^{-k+1} = 2^{1-k} \text{Tr} \mathbf{D}$$

so that all the diagonal elements are identically equal to the average of the diagonal elements of  $D$ .  $\square$

### C. Proofs of Theorems in Section 3

*Proof of Theorem 3.3.* Let  $UV^\top = W\Sigma Y^\top$  be a compact SVD of  $UV^\top$ . Define  $\tilde{U} := W\Sigma^{1/2}$  and  $\tilde{V} := Y\Sigma^{1/2}$  and observe that  $\tilde{U}\tilde{V}^\top = UV^\top$ . Furthermore, let  $G_{\tilde{U}} = \tilde{U}^\top\tilde{U}$  and  $G_{\tilde{V}} = \tilde{V}^\top\tilde{V}$  be their Gram matrices. Observe that  $G_{\tilde{U}} = G_{\tilde{V}} = \Sigma$ . Hence, by Theorem 2.2, there exists a rotation  $Q$  such that  $\tilde{V} := \tilde{V}Q$  and  $\tilde{U} := \tilde{U}Q$  are equalized, with  $\|\tilde{u}_i\|^2 = \|\tilde{v}_i\|^2 = \frac{1}{r} \text{Tr } \Sigma$ .  $\square$

*Proof of Theorem 3.4.* Define

$$\mathbf{n}_{\mathbf{u},\mathbf{v}} = (\|\mathbf{u}_1\|\|\mathbf{v}_1\|, \dots, \|\mathbf{u}_r\|\|\mathbf{v}_r\|)$$

and observe that

$$\begin{aligned} R(\mathbf{U}, \mathbf{V}) &= \lambda \sum_{i=1}^r \|\mathbf{u}_i\|^2 \|\mathbf{v}_i\|^2 \\ &= \frac{\lambda}{r} \|\mathbf{n}_{\mathbf{u},\mathbf{v}}\|^2 \|\mathbf{1}_r\|^2 \geq \frac{\lambda}{r} \langle \mathbf{n}_{\mathbf{u},\mathbf{v}}, \mathbf{1}_r \rangle^2 \\ &= \frac{\lambda}{r} \left( \sum_{i=1}^r \|\mathbf{u}_i\|\|\mathbf{v}_i\| \right)^2 \end{aligned}$$

where the inequality is due to Cauchy-Schwartz, and it holds with equality if and only if  $\mathbf{n}_{\mathbf{u},\mathbf{v}}$  is parallel to  $\mathbf{1}_r$ . Let  $(\tilde{\mathbf{U}}, \tilde{\mathbf{V}})$  be a global optima of Problem 6. The inequality above together with Theorem 3.3 imply that  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  should be jointly equalized up to dilation transformations, hence the first equality claimed by the theorem.

To see the second equality, note that if  $\mathbf{U}$  and  $\mathbf{V}$  are jointly equalized, then

$$\|\mathbf{u}_i\|^2 = \|\mathbf{v}_i\|^2 = \frac{1}{r} \text{Tr } \Sigma,$$

where  $\Sigma$  is the matrix of singular values of  $UV^\top$ . Hence,

$$\begin{aligned} R(\mathbf{U}, \mathbf{V}) &= \frac{\lambda}{r} \left( \sum_{i=1}^r \|\mathbf{u}_i\|\|\mathbf{v}_i\| \right)^2 = \frac{\lambda}{r} \left( \frac{1}{r} \sum_{i=1}^r \text{Tr } \Sigma \right)^2 \\ &= \frac{\lambda}{r} (\text{Tr } \Sigma)^2 \end{aligned}$$

which is equal to  $\frac{\lambda}{r} \|\tilde{\mathbf{U}}\tilde{\mathbf{V}}^\top\|_*^2$  as claimed.  $\square$

*Proof of Theorem 3.6.* By Theorem 3.4, if  $(\mathbf{X}, \mathbf{Y})$  is an optimum of Problem 6, then it holds that

$$\lambda \sum_{i=1}^r \|\mathbf{x}_i\|^2 \|\mathbf{y}_i\|^2 = \frac{\lambda}{r} \|\mathbf{X}\mathbf{Y}^\top\|_*^2.$$

Hence, Problem 6 reduces to the following problem:

$$\min_{\mathbf{X} \in \mathbb{R}^{d_1 \times r}, \mathbf{Y} \in \mathbb{R}^{d_2 \times r}} \|\mathbf{M} - \mathbf{X}\mathbf{Y}^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{X}\mathbf{Y}^\top\|_*^2 \quad (14)$$

Let  $\mathbf{M} = \mathbf{U}_M \Sigma_M \mathbf{V}_M^\top$  and  $\mathbf{W} := \mathbf{X}\mathbf{Y}^\top = \mathbf{U}_W \Sigma_W \mathbf{V}_W^\top$  be full SVDs of  $\mathbf{M}$  and  $\mathbf{W}$  respectively, such that  $\sigma_i(\mathbf{M}) \geq \sigma_{i+1}(\mathbf{M})$  and  $\sigma_i(\mathbf{W}) \geq \sigma_{i+1}(\mathbf{W})$  for all  $i \in [d-1]$  where  $d = \min\{d_1, d_2\}$ . Rewriting objective of Problem 14 in terms of these decompositions,

$$\begin{aligned} &\|\mathbf{M} - \mathbf{X}\mathbf{Y}^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{X}\mathbf{Y}^\top\|_*^2 \\ &= \|\mathbf{U}_M \Sigma_M \mathbf{V}_M^\top - \mathbf{U}_W \Sigma_W \mathbf{V}_W^\top\|_F^2 + \frac{\lambda}{r} \|\mathbf{U}_W \Sigma_W \mathbf{V}_W^\top\|_*^2 \\ &= \|\Sigma_M - \mathbf{U}' \Sigma_W \mathbf{V}'^\top\|_F^2 + \frac{\lambda}{r} \|\Sigma_W\|_*^2 \\ &= \|\Sigma_M\|_F^2 + \|\Sigma_W\|_F^2 - 2\langle \Sigma_M, \mathbf{U}' \Sigma_W \mathbf{V}'^\top \rangle + \frac{\lambda}{r} \|\Sigma_W\|_*^2 \end{aligned}$$

where  $\mathbf{U}' = \mathbf{U}_M^\top \mathbf{U}_W$ . By Von Neumann's trace inequality, for a fixed  $\Sigma_W$  we have that  $\langle \Sigma_M, \mathbf{U}' \Sigma_W \mathbf{V}'^\top \rangle \leq \sum_{i=1}^d \sigma_i(\mathbf{M}) \sigma_i(\mathbf{W})$ , where the equality is achieved when  $\Sigma_i(\mathbf{W})$  have the same ordering as  $\Sigma_i(\mathbf{M})$  and  $\mathbf{U}' = \mathbf{I}$ , i.e.  $\mathbf{U}_M = \mathbf{U}_W$ . Now, Problem 14 is reduced to

$$\begin{aligned} &\min_{\substack{\|\Sigma_W\|_0 \leq r, \\ \Sigma_W \geq 0}} \|\Sigma_M - \Sigma_W\|_F^2 + \frac{\lambda}{r} \|\Sigma_W\|_*^2 \\ &= \min_{\sigma \in \mathbb{R}_+^r} \sum_{i=1}^r (\sigma_i(\mathbf{M}) - \bar{\sigma}_i)^2 + \sum_{i=r+1}^d \sigma_i^2(\mathbf{M}) + \frac{\lambda}{r} \left( \sum_{i=1}^r \bar{\sigma}_i \right)^2 \end{aligned}$$

The Lagrangian is given by

$$\begin{aligned} L(\bar{\lambda}, \alpha) &= \sum_{i=1}^r (\sigma_i(\mathbf{M}) - \bar{\sigma}_i)^2 + \sum_{i=r+1}^d \sigma_i^2(\mathbf{M}) \\ &\quad + \frac{\lambda}{r} \left( \sum_{i=1}^r \bar{\sigma}_i \right)^2 - \sum_{i=1}^r \alpha_i \bar{\sigma}_i \end{aligned}$$

The KKT conditions ensures that  $\forall i = 1, \dots, r$ ,

$$\begin{aligned} &\bar{\sigma}_i \geq 0, \alpha_i \geq 0, \bar{\sigma}_i \alpha_i = 0 \\ &2(\bar{\sigma}_i - \sigma_i(\mathbf{M})) + \frac{2\lambda}{r} \left( \sum_{i=1}^r \bar{\sigma}_i \right) - \alpha_i = 0 \end{aligned}$$

Let  $\rho = |\{i : \bar{\sigma}_i > 0\}| \leq r$  be the number of nonzero  $\bar{\sigma}_i$ . For

$i = 1, \dots, \rho$  we have  $\alpha_i = 0$ , hence

$$\begin{aligned} \bar{\sigma}_i + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \bar{\sigma}_i \right) &= \sigma_i(\mathbf{M}) \\ \implies (\mathbf{I}_{\rho} + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^{\top}) \bar{\sigma}_{1:\rho} &= \sigma_{1:\rho}(\mathbf{M}) \\ \implies \bar{\sigma}_{1:\rho} &= (\mathbf{I}_{\rho} - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^{\top}) \sigma_{1:\rho}(\mathbf{M}) \\ \implies \bar{\sigma}_{1:\rho} &= \sigma_{1:\rho}(\mathbf{M}) - \frac{\lambda\rho\kappa_{\rho}}{r + \lambda\rho} \mathbf{1}_{\rho} \\ \implies \Sigma_{\mathbf{W}} &= (\Sigma_{\mathbf{M}} - \frac{\lambda\rho\kappa_{\rho}}{r + \lambda\rho} \mathbf{I}_d)_+ \end{aligned}$$

where  $\kappa_{\rho} = \frac{1}{\rho} \sum_{i=1}^{\rho} \sigma_i(\mathbf{M})$  and the second implication holds since  $(\mathbf{I}_{\rho} + \frac{\lambda}{r} \mathbf{1}\mathbf{1}^{\top})^{-1} = \mathbf{I}_{\rho} - \frac{\lambda}{r + \lambda\rho} \mathbf{1}\mathbf{1}^{\top}$ . It only remains to find the optimal  $\rho$ . Let's define the function

$$\begin{aligned} g(\rho) &:= \sum_{i=1}^{\rho} (\sigma_i(\mathbf{M}) - \bar{\sigma}_i)^2 + \sum_{i=\rho+1}^d \sigma_i^2(\mathbf{M}) + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \bar{\sigma}_i \right)^2 \\ &= \sum_{i=1}^{\rho} \left( \frac{\lambda \sum_{k=1}^{\rho} \sigma_k(\mathbf{M})}{r + \lambda\rho} \right)^2 + \sum_{i=\rho+1}^d \sigma_i(\mathbf{M})^2 \\ &\quad + \frac{\lambda}{r} \left( \sum_{i=1}^{\rho} \left( \sigma_i(\mathbf{M}) - \frac{\lambda \sum_{k=1}^{\rho} \sigma_k(\mathbf{M})}{r + \lambda\rho} \right) \right)^2. \end{aligned}$$

By Lemma A.3,  $g(\rho)$  is monotonically non-increasing in  $\rho$ , hence  $\rho$  should be the largest *feasible* integer, i.e.

$$\rho = \max\{j : \sigma_j > \frac{\lambda j \kappa_j}{r + \lambda j}\}.$$

□

## D. Proofs of Theorems in Sections 4

In this section for ease of notation we let  $\lambda_i$  denote  $\lambda_i(\mathbf{M})$ . Furthermore, with slight abuse of notation we let  $f(\mathbf{U})$ ,  $\ell(\mathbf{U})$  and  $R(\mathbf{U})$  denote the objective, the loss function and the regularizer, respectively.

It is easy to see that the gradient of the objective of Problem 4 is given by

$$\nabla f(\mathbf{U}) = 4(\mathbf{U}\mathbf{U}^{\top} - \mathbf{M})\mathbf{U} + 4\lambda\mathbf{U} \operatorname{diag}(\mathbf{U}^{\top}\mathbf{U}).$$

We first make the following important observation about the critical points of Problem 4.

**Lemma D.1.** If  $\mathbf{U}$  is a critical point of Problem 4, then it holds that  $\mathbf{U}\mathbf{U}^{\top} \preceq \mathbf{M}$ .

*Proof of Lemma D.1.* Since  $\nabla f(\mathbf{U}) = 0$ , we have that

$$(\mathbf{M} - \mathbf{U}\mathbf{U}^{\top})\mathbf{U} = \lambda\mathbf{U} \operatorname{diag}(\mathbf{U}^{\top}\mathbf{U})$$

multiply both sides from right by  $\mathbf{U}^{\top}$  and rearrange to get

$$\mathbf{M}\mathbf{U}\mathbf{U}^{\top} = \mathbf{U}\mathbf{U}^{\top}\mathbf{U}\mathbf{U}^{\top} + \lambda\mathbf{U} \operatorname{diag}(\mathbf{U}^{\top}\mathbf{U})\mathbf{U}^{\top} \quad (15)$$

Note that the right hand side is symmetric, which implies that the left hand side must be symmetric as well, i.e.

$$\mathbf{M}\mathbf{U}\mathbf{U}^{\top} = (\mathbf{M}\mathbf{U}\mathbf{U}^{\top})^{\top} = \mathbf{U}\mathbf{U}^{\top}\mathbf{M},$$

so that  $\mathbf{M}$  and  $\mathbf{U}\mathbf{U}^{\top}$  commute. Note that in Equation (15),  $\mathbf{U} \operatorname{diag}(\mathbf{U}^{\top}\mathbf{U})\mathbf{U}^{\top} \succeq 0$ . Thus,  $\mathbf{M}\mathbf{U}\mathbf{U}^{\top} \succeq \mathbf{U}\mathbf{U}^{\top}\mathbf{U}\mathbf{U}^{\top}$ . Let  $\mathbf{U}\mathbf{U}^{\top} = \mathbf{W}\mathbf{\Gamma}\mathbf{W}^{\top}$  be a compact eigendecomposition of  $\mathbf{U}\mathbf{U}^{\top}$ . We get

$$\mathbf{M}\mathbf{U}\mathbf{U}^{\top} = \mathbf{M}\mathbf{W}\mathbf{\Gamma}\mathbf{W}^{\top} \succeq \mathbf{U}\mathbf{U}^{\top}\mathbf{U}\mathbf{U}^{\top} = \mathbf{W}\mathbf{\Gamma}^2\mathbf{W}^{\top}.$$

Multiplying from right and left by  $\mathbf{W}\mathbf{\Gamma}^{-1}$  and  $\mathbf{W}^{\top}$  respectively, we have that

$$\mathbf{W}^{\top}\mathbf{M}\mathbf{W} \succeq \mathbf{\Gamma}$$

which completes the proof. □

Lemma D.1 allows us to bound different norms of the critical points, as will be seen later in the proofs.

To explore the landscape properties of Problem 4, we first focus on the non-equalized critical points in Lemma D.2. We show that the set of non-equalized critical points does not include any local optima. Furthermore, all such points are strict saddles. Therefore, we turn our focus to the equalized critical points in Lemma D.3. We show all such points inherit the eigenspace of the input matrix  $\mathbf{M}$ . This allows us to give a closed-form characterization of all the equalized critical points in terms of the eigendecomposition of  $\mathbf{M}$ . We then show that if  $\lambda$  is chosen appropriately, all such critical points that are not global optima, are strict saddle points.

**Lemma D.2.** All local minima of Problem 4 are equalized. Moreover, all critical points that are not equalized, are strict saddle points.

*Proof of Lemma D.2.* We show that if  $\mathbf{U}$  is not equalized, then any  $\epsilon$ -neighborhood of  $\mathbf{U}$  contains a point with objective strictly smaller than  $f(\mathbf{U})$ . More formally, for any  $\epsilon > 0$ , we exhibit a rotation  $\mathbf{Q}_{\epsilon}$  such that  $\|\mathbf{U} - \mathbf{U}\mathbf{Q}_{\epsilon}\|_F \leq \epsilon$  and  $f(\mathbf{U}\mathbf{Q}_{\epsilon}) < f(\mathbf{U})$ . Let  $\mathbf{U}$  be a critical point of Problem 4 that is not equalized, i.e. there exists two columns of  $\mathbf{U}$  with different norms. Without loss of generality, let  $\|\mathbf{u}_1\| > \|\mathbf{u}_2\|$ . We design a rotation matrix  $\mathbf{Q}$  such that it is almost an isometry, but it moves mass from  $\mathbf{u}_1$  to  $\mathbf{u}_2$ . Consequently, the new factor becomes “less un-equalized” and achieves a smaller regularizer, while preserving the value of the loss. To that end, define

$$\mathbf{Q}_{\delta} := \begin{bmatrix} \sqrt{1 - \delta^2} & -\delta & 0 \\ \delta & \sqrt{1 - \delta^2} & 0 \\ 0 & 0 & \mathbf{I}_{r-2} \end{bmatrix}$$

and let  $\hat{U} := \mathbf{U}\mathbf{Q}_\delta$ . It is easy to verify that  $\mathbf{Q}_\epsilon$  is indeed a rotation. First, we show that for any  $\epsilon$ , as long as  $\delta^2 \leq \frac{\epsilon^2}{2\text{Tr}(\mathbf{M})}$ , we have  $\hat{U} \in \mathcal{B}_\epsilon(\mathbf{U})$ :

$$\begin{aligned} \|\mathbf{U} - \hat{U}\|_F^2 &= \sum_{i=1}^r \|\mathbf{u}_i - \hat{\mathbf{u}}_i\|^2 \\ &= \|\mathbf{u}_1 - \sqrt{1-\delta^2}\mathbf{u}_1 - \delta\mathbf{u}_2\|^2 \\ &\quad + \|\mathbf{u}_2 - \sqrt{1-\delta^2}\mathbf{u}_2 + \delta\mathbf{u}_1\|^2 \\ &= 2(1 - \sqrt{1-\delta^2})(\|\mathbf{u}_1\|^2 + \|\mathbf{u}_2\|^2) \\ &\leq 2\delta^2 \text{Tr}(\mathbf{M}) \leq \epsilon^2 \end{aligned}$$

where the second to last inequality follows from Lemma D.1, because  $\|\mathbf{u}_1\|^2 + \|\mathbf{u}_2\|^2 \leq \|\mathbf{U}\|_F^2 = \text{Tr}(\mathbf{U}\mathbf{U}^\top) \leq \text{Tr}(\mathbf{M})$ , and also the fact that  $1 - \sqrt{1-\delta^2} = \frac{1-\delta^2}{1+\sqrt{1-\delta^2}} \leq \delta^2$ .

Next, we show that for small enough  $\delta$ , the value of the function at  $\hat{U}$  is strictly smaller than that of  $\mathbf{U}$ . Observe that

$$\begin{aligned} \|\hat{\mathbf{u}}_1\|^2 &= (1-\delta^2)\|\mathbf{u}_1\|^2 + \delta^2\|\mathbf{u}_2\|^2 + 2\delta\sqrt{1-\delta^2}\mathbf{u}_1^\top\mathbf{u}_2 \\ \|\hat{\mathbf{u}}_2\|^2 &= (1-\delta^2)\|\mathbf{u}_2\|^2 + \delta^2\|\mathbf{u}_1\|^2 - 2\delta\sqrt{1-\delta^2}\mathbf{u}_1^\top\mathbf{u}_2 \end{aligned}$$

and the remaining columns will not change, i.e. for  $i = 3, \dots, r$ ,  $\hat{\mathbf{u}}_i = \mathbf{u}_i$ . Together with the fact that  $\mathbf{Q}_\delta$  preserves the norms, i.e.  $\|\mathbf{U}\|_F = \|\mathbf{U}\mathbf{Q}_\delta\|_F$ , we get

$$\|\hat{\mathbf{u}}_1\|^2 + \|\hat{\mathbf{u}}_2\|^2 = \|\mathbf{u}_1\|^2 + \|\mathbf{u}_2\|^2. \quad (16)$$

Let  $\delta = -c \cdot \text{sgn}(\mathbf{u}_1^\top\mathbf{u}_2)$  for a small enough  $c > 0$  such that  $\|\mathbf{u}_2\| < \|\hat{\mathbf{u}}_2\| \leq \|\hat{\mathbf{u}}_1\| < \|\mathbf{u}_1\|$ . Using Equation (16), This implies that  $\|\hat{\mathbf{u}}_1\|^4 + \|\hat{\mathbf{u}}_2\|^4 < \|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4$ , which in turn gives us  $R(\hat{U}) < R(\mathbf{U})$  and hence  $f(\hat{U}) < f(\mathbf{U})$ . Therefore, a non-equalized critical point cannot be local minimum, hence the first claim of the lemma.

We now prove the second part of the lemma. Let  $\mathbf{U}$  be a critical point that is not equalized. To show that  $\mathbf{U}$  is a strict saddle point, it suffices to show that the Hessian has a negative eigenvalue. In here, we exhibit a curve along which the second directional derivative is negative. Assume, without loss of generality that  $\|\mathbf{u}_1\| > \|\mathbf{u}_2\|$  and consider the curve

$$\Delta(t) := [(\sqrt{1-t^2}-1)\mathbf{u}_1 + t\mathbf{u}_2, (\sqrt{1-t^2}-1)\mathbf{u}_2 - t\mathbf{u}_1, \mathbf{0}_{d,r-2}]$$

It is easy to check that for any  $t \in \mathbb{R}$ ,  $\ell(\mathbf{U} + \Delta(t)) = \ell(\mathbf{U})$  since  $\mathbf{U} + \Delta(t)$  is essentially a rotation on  $\mathbf{U}$  and  $\ell$  is invariant under rotations. Observe that

$$\begin{aligned} g(t) &:= f(\mathbf{U} + \Delta(t)) \\ &= f(\mathbf{U}) + \|\sqrt{1-t^2}\mathbf{u}_1 + t\mathbf{u}_2\|^4 - \|\mathbf{u}_1\|^4 \\ &\quad + \|\sqrt{1-t^2}\mathbf{u}_2 - t\mathbf{u}_1\|^4 - \|\mathbf{u}_2\|^4 \\ &= f(\mathbf{U}) - 2t^2(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 8t^2(\mathbf{u}_1\mathbf{u}_2)^2 \\ &\quad + 4t^2\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 + 4t\sqrt{1-t^2}\mathbf{u}_1^\top\mathbf{u}_2(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) + O(t^3). \end{aligned}$$

The derivative of  $g$  then is given as

$$\begin{aligned} g'(t) &= -4t(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 16t(\mathbf{u}_1\mathbf{u}_2)^2 + 8t\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 \\ &\quad + 4(\sqrt{1-t^2} - \frac{t^2}{\sqrt{1-t^2}})(\mathbf{u}_1^\top\mathbf{u}_2)(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) + O(t^2). \end{aligned}$$

Since  $\mathbf{U}$  is a critical point and  $f$  is continuously differentiable, it should hold that  $g'(0) = 4(\mathbf{u}_1^\top\mathbf{u}_2)(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2) = 0$ . Since by assumption  $\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2 > 0$ , it should be the case that  $\mathbf{u}_1^\top\mathbf{u}_2 = 0$ . We now consider the second order directional derivative:

$$\begin{aligned} g''(0) &= -4(\|\mathbf{u}_1\|^4 + \|\mathbf{u}_2\|^4) + 16(\mathbf{u}_1\mathbf{u}_2)^2 + 8\|\mathbf{u}_1\|^2\|\mathbf{u}_2\|^2 \\ &= -4(\|\mathbf{u}_1\|^2 - \|\mathbf{u}_2\|^2)^2 < 0 \end{aligned}$$

which completes the proof.  $\square$

We now focus on the critical points that are equalized, i.e. points  $\mathbf{U}$  such that  $\nabla f(\mathbf{U}) = 0$  and  $\text{diag}(\mathbf{U}^\top\mathbf{U}) = \frac{\|\mathbf{U}\|_F^2}{r}\mathbf{I}$ .

**Lemma D.3.** Assume that  $\lambda < \frac{r\lambda_r}{\sum_{i=1}^r \lambda_i - r\lambda_r}$ . Then all equalized local minima are global. All other equalized critical points are strict saddle points.

*Proof of Lemma D.3.* Let  $\mathbf{U} = \mathbf{W}\Sigma\mathbf{V}^\top$  be a compact SVD of the rank- $r'$  weight matrix  $\mathbf{U}$ . We have:

$$\begin{aligned} \nabla f(\mathbf{U}) &= 4(\mathbf{U}\mathbf{U}^\top - \mathbf{M})\mathbf{U} + 4\lambda\mathbf{U} \text{diag}(\mathbf{U}^\top\mathbf{U}) = 0 \\ \implies \mathbf{U}\mathbf{U}^\top\mathbf{U} + \frac{\lambda\|\mathbf{U}\|_F^2}{r}\mathbf{U} &= \mathbf{M}\mathbf{U} \\ \implies \mathbf{W}\Sigma^3\mathbf{V}^\top + \frac{\lambda\|\Sigma\|_F^2}{r}\mathbf{W}\Sigma\mathbf{V}^\top &= \mathbf{M}\mathbf{W}\Sigma\mathbf{V}^\top \\ \implies \Sigma^2 + \frac{\lambda\|\Sigma\|_F^2}{r}\mathbf{I} &= \mathbf{W}^\top\mathbf{M}\mathbf{W} \end{aligned}$$

Since the left hand side of the above equality is diagonal, it implies that  $\mathbf{W} \in \mathbb{R}^{d \times r'}$  corresponds to some  $r'$  eigenvectors of  $\mathbf{M}$ . Let  $\mathcal{E} \subseteq [d]$ ,  $|\mathcal{E}| = r'$  denote the set of eigenvectors of  $\mathbf{M}$  that are present in  $\mathbf{W}$ . Note that the above is equivalent of the following system of linear equations:

$$(\mathbf{I} + \frac{\lambda}{r}\mathbf{1}\mathbf{1}^\top)\sigma^2 = \vec{\lambda},$$

where  $\sigma^2 := \text{diag}(\Sigma^2)$  and  $\vec{\lambda} = \text{diag}(\mathbf{W}^\top\mathbf{M}\mathbf{W})$ . By Lemma A.2, the solution to this linear system is given by

$$\sigma^2 = (\mathbf{I} - \frac{\lambda}{r + \lambda r'})\vec{\lambda}. \quad (17)$$

The set  $\mathcal{E}$  belongs to one of the following categories:

1.  $\mathcal{E} = [r']$ ,  $r' = \rho$
2.  $\mathcal{E} = [r']$ ,  $r' < \rho$
3.  $\mathcal{E} \neq [r']$

The case  $\mathcal{E} = [r']$ ,  $r' > \rho$  is excluded from the above partition, since whenever  $\mathcal{E} = [r']$ , it should hold that  $r' \leq \rho$ . To see this, note that due to  $\mathbf{U} = \mathbf{W}\Sigma\mathbf{V}^\top$  being a compact SVD of  $\mathbf{M}$ , it holds that  $\sigma_j > 0$  for all  $j \in [r']$ . Specifically for  $j = r'$ , plugging  $\sigma_{r'} > 0$  back to Equation (17), we get

$$\lambda_{r'} > \frac{\lambda \sum_{i=1}^{r'} \lambda_i}{r + \lambda r'} = \frac{\lambda r' \kappa_{r'}}{r + \lambda r'}.$$

Then it follows from definition of  $\rho$  in Theorem 2.4 that  $r' \leq \rho$ . We provide a case by case analysis for the above partition here.

**Case 1.** [ $\mathcal{E} = [r']$ ,  $r' = \rho$ ] When  $\mathbf{W}$  corresponds to the top- $\rho$  eigenvectors of  $\mathbf{M}$ , we retrieve the global optimal solution described by Theorem 2.4. Therefore, all such critical points are global minima.

**Case 2.** [ $\mathcal{E} = [r']$ ,  $r' < \rho$ ] Let  $\mathbf{W}_r := [\mathbf{W}, \mathbf{W}_\perp]$  be the top- $r$  eigenvectors of  $\mathbf{M}$  and  $\mathbf{V}_\perp$  span the orthogonal subspace of  $\mathbf{V}$ , i.e.  $\mathbf{V}_r := [\mathbf{V}, \mathbf{V}_\perp]$  be an orthonormal basis for  $\mathbb{R}^r$ . Define  $\mathbf{U}(t) = \mathbf{W}_r \Sigma' \mathbf{V}_r^\top$  where  $\sigma'_i = \sqrt{\sigma_i^2 + t^2}$  for  $i \leq r$ . Observe that

$$\mathbf{U}(t)^\top \mathbf{U}(t) = \mathbf{V} \Sigma \mathbf{V}^\top + t^2 \mathbf{V}_r^\top \mathbf{V}_r = \mathbf{U}^\top \mathbf{U} + t^2 \mathbf{I}_r$$

so that for all  $t$ , the parametric curve  $\mathbf{U}(t)$  is equalized. The value of the loss function at  $\mathbf{U}(t)$  is given by:

$$\begin{aligned} \ell(\mathbf{U}(t)) &= \sum_{i=1}^r (\lambda_i - \sigma_i^2 - t^2)^2 + \sum_{i=r+1}^d (\lambda_i)^2 \\ &= \ell(\mathbf{U}) + r t^4 - 2t^2 \sum_{i=1}^r (\lambda_i - \sigma_i^2). \end{aligned}$$

Furthermore, since  $\mathbf{U}(t)$  is equalized, we obtain the following form for the regularizer:

$$\begin{aligned} R(\mathbf{U}(t)) &= \frac{\lambda}{r} \|\mathbf{U}(t)\|_F^4 = \frac{\lambda}{r} (\|\mathbf{U}\|_F^2 + r t^2)^2 \\ &= \ell(\mathbf{U}) + \lambda r t^4 + 2\lambda t^2 \|\mathbf{U}\|_F^2. \end{aligned}$$

Now define  $g(t) := \ell(\mathbf{U}(t)) + R(\mathbf{U}(t))$  and observe

$$\begin{aligned} g(t) &= \ell(\mathbf{U}) + R(\mathbf{U}) + r t^4 - 2t^2 \sum_{i=1}^r (\lambda_i - \sigma_i^2) \\ &\quad + \lambda r t^4 + 2\lambda t^2 \|\mathbf{U}\|_F^2. \end{aligned}$$

It is easy to verify that  $g'(0) = 0$ . Moreover, the second derivative of  $g$  at the origin is given as:

$$\begin{aligned} g''(0) &= -4 \sum_{i=1}^r (\lambda_i - \sigma_i^2) + 4\lambda \|\mathbf{U}\|_F^2 \\ &= -4 \sum_{i=1}^r \lambda_i + 4(1 + \lambda) \|\mathbf{U}\|_F^2 \\ &= -4 \sum_{i=1}^r \lambda_i + 4 \frac{r + r\lambda}{r + \lambda r'} \sum_{i=1}^{r'} \lambda_i \end{aligned}$$

where the last equality follows from the fact Equation (17) and the fact that  $\|\mathbf{U}\|_F^2 = \sum_{i=1}^{r'} \sigma_i^2$ . To get a sufficient condition for  $\mathbf{U}$  to be a strict saddle point, we set  $g''(0) < 0$ :

$$\begin{aligned} &-4 \sum_{i=r'+1}^r \lambda_i + 4 \frac{(r - r')\lambda}{r + \lambda r'} \sum_{i=1}^{r'} \lambda_i < 0 \\ \implies &\frac{(r - r')\lambda}{r + \lambda r'} \sum_{i=1}^{r'} \lambda_i < \sum_{i=r'+1}^r \lambda_i \\ \implies &\lambda < \frac{(r + \lambda r') \sum_{i=r'+1}^r \lambda_i}{(r - r') \sum_{i=1}^{r'} \lambda_i} \\ \implies &\lambda \left(1 - \frac{r' \sum_{i=r'+1}^r \lambda_i}{(r - r') \sum_{i=1}^{r'} \lambda_i}\right) < \frac{r \sum_{i=r'+1}^r \lambda_i}{(r - r') \sum_{i=1}^{r'} \lambda_i} \\ \implies &\lambda < \frac{r \sum_{i=r'+1}^r \lambda_i}{(r - r') \sum_{i=1}^{r'} \lambda_i - r' \sum_{i=r'+1}^r \lambda_i} \\ \implies &\lambda < \frac{r h(r')}{\sum_{i=1}^{r'} (\lambda_i - h(r'))} \end{aligned}$$

where  $h(r') := \frac{\sum_{i=r'+1}^r \lambda_i}{r - r'}$  is the average of the eigenvalues  $\lambda_{r'+1}, \dots, \lambda_r$ . It is easy to see that the right hand side is monotonically decreasing with  $r'$ , since  $h(r')$  monotonically decrease with  $r'$ . Hence, it suffices to make sure that  $\lambda$  is smaller than the right hand side for the choice of  $r' = r - 1$ , i.e.  $\lambda < \frac{r\lambda_r}{\sum_{i=1}^r (\lambda_i - \lambda_r)}$ .

**Case 3.** [ $\mathcal{E} \neq [r']$ ] We show that all such critical points are strict saddle points. Let  $\mathbf{w}'$  be one of the top  $r'$  eigenvectors that are missing in  $\mathbf{W}$ . Let  $j \in \mathcal{E}$  be such that  $\mathbf{w}_j$  is not among the top  $r'$  eigenvectors of  $\mathbf{M}$ . For any  $t \in [0, 1]$ , let  $\mathbf{W}(t)$  be identical to  $\mathbf{W}$  in all the columns but the  $j^{\text{th}}$  one, where  $\mathbf{w}_j(t) = \sqrt{1 - t^2} \mathbf{w}_j + t \mathbf{w}'$ . Note that  $\mathbf{W}(t)$  is still an orthogonal matrix for all values of  $t$ . Define the parametrized curve  $\mathbf{U}(t) := \mathbf{W}(t) \Sigma \mathbf{V}^\top$  for  $t \in [0, 1]$  and observe that:

$$\begin{aligned} \|\mathbf{U} - \mathbf{U}(t)\|_F^2 &= \sigma_j^2 \|\mathbf{w}_j - \mathbf{w}_j(t)\|^2 \\ &= 2\sigma_j^2 (1 - \sqrt{1 - t^2}) \leq t^2 \text{Tr } \mathbf{M} \end{aligned}$$

That is, for any  $\epsilon > 0$ , there exist a  $t > 0$  such that  $\mathbf{U}(t)$  belongs to the  $\epsilon$ -ball around  $\mathbf{U}$ . We show that  $f(\mathbf{U}(t))$  is strictly smaller than  $f(\mathbf{U})$ , which means  $\mathbf{U}$  cannot be a local minimum. Note that this construction of  $\mathbf{U}(t)$  guarantees that  $R(\mathbf{U}') = R(\mathbf{U})$ . In particular, it is easy to see that  $\mathbf{U}(t)^\top \mathbf{U}(t) = \mathbf{U}^\top \mathbf{U}$ , so that  $\mathbf{U}(t)$  remains equalized for all values of  $t$ . Moreover, we have that

$$\begin{aligned} f(\mathbf{U}(t)) - f(\mathbf{U}) &= \|\mathbf{M} - \mathbf{U}(t)\mathbf{U}(t)^\top\|_F^2 - \|\mathbf{M} - \mathbf{U}\mathbf{U}^\top\|_F^2 \\ &= -2 \text{Tr}(\Sigma^2 \mathbf{W}(t)^\top \mathbf{M} \mathbf{W}(t)) + 2 \text{Tr}(\Sigma^2 \mathbf{W}^\top \mathbf{M} \mathbf{W}) \\ &= -2\sigma_j^2 t^2 (\mathbf{w}_j(t)^\top \mathbf{M} \mathbf{w}_j(t) - \mathbf{w}_j^\top \mathbf{M} \mathbf{w}_j) < 0, \end{aligned}$$

where the last inequality follows because by construction  $\mathbf{w}_j(t)^\top \mathbf{M} \mathbf{w}_j(t) > \mathbf{w}_j^\top \mathbf{M} \mathbf{w}_j$ . Define  $g(t) := f(\mathbf{U}(t)) = \ell(\mathbf{U}(t)) + R(\mathbf{U}(t))$ . To see that such saddle points are non-degenerate, it suffices to show  $g''(0) < 0$ . It is easy to check that the second directional derivative at the origin is given by

$$g''(0) = -4\sigma_j^2(\mathbf{w}_j(t)^\top \mathbf{M} \mathbf{w}_j(t) - \mathbf{w}_j^\top \mathbf{M} \mathbf{w}_j) < 0,$$

which completes the proof. □

*Proof of Lemma 4.1.* Follows from Lemma D.2 □

*Proof of Theorem 4.3.* Follows from Lemma D.2 and Lemma D.3. □