

A Numerical Simulations

In this section we present simulation results for performance of gradient descent over $f(\mathbf{U})$. We consider measurements $y_i = \langle \mathbf{A}_i, \mathbf{X}^* \rangle$, where \mathbf{A}_i are i.i.d Gaussian with each entry distributed as $\mathcal{N}(0, 1/m)$. \mathbf{X}^* is a 100×100 rank r random p.s.d matrix with $\|\mathbf{X}^*\|_F = 1$. r is varied from 1 to 20 in the experiments.

We consider both standard gradient descent and noisy gradient descent (4) with step size $\frac{1}{\|\mathbf{U}\|_2}$. We add noise of magnitude $1e-4$ for the noisy gradient updates. Each method is run until convergence (max of 200 iterations). Let the output of gradient descent be $\hat{\mathbf{U}}$. A run of this experiment is considered success if the final error $\|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{X}^*\|_F \leq 1e-2$. Each experiment is repeated for 20 times and average probability of success is computed.

We repeat the above procedure starting from both random initialization and SVD initialization. For SVD initialization, the initial point is set to be the rank r approximation of $\sum_{i=1}^m y_i \mathbf{A}_i$ as suggested by Jain et al. [15]. In figure 2 we have the plots for the cases discussed above. All of them have phase transition around number of samples $m = 2 \cdot n \cdot r$. This is in agreement with the results in Section 3. $f(\mathbf{U})$ has no local minima once $m \geq 2 \cdot n \cdot r$ and random initialization has same performance as SVD initialization.

In figure 3, the left two plots show error $\|\hat{\mathbf{U}}\hat{\mathbf{U}}^\top - \mathbf{X}^*\|_F / \|\mathbf{X}^*\|_F$ behaves with varying rank and number of samples for random and SVD initializations. The rightmost plot shows the phase transition for rank 10 case for all the methods. Again we notice no significant difference between these methods.

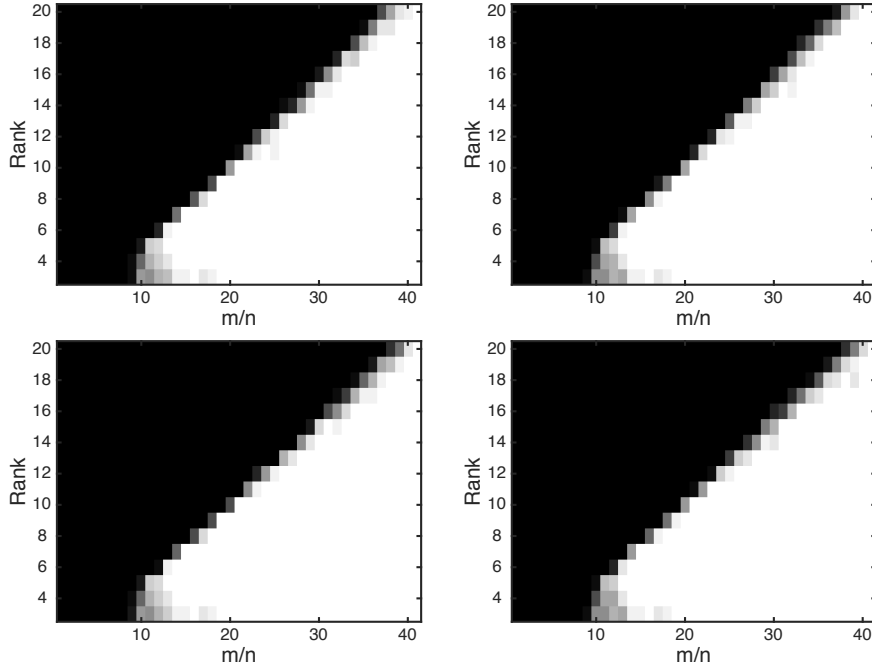


Figure 2: This figure plots the success probability for increasing number of samples m and various values of rank r . The plots on the top are for gradient descent, left for random initialization and the right for SVD initialization. Similarly the bottom plots are for the noisy gradient descent. We notice no significant difference between all these settings. They all have phase transition around $m = 2 \cdot n \cdot r$.

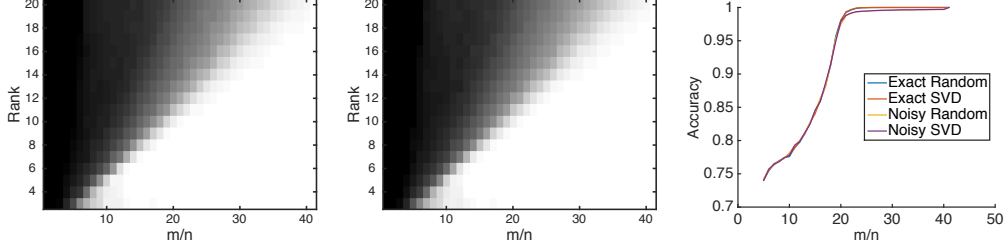


Figure 3: This figure plots the error $\|\hat{U}\hat{U}^\top - X^*\|_F / \|X^*\|_F$ for increasing number of samples m . The left plot is for gradient descent with random initialization, center plot corresponds to SVD initialization. Again we notice no difference in error for these two settings. The rightmost figure shows phase transition of low rank recovery for all the different settings when X^* is rank 10.

B Proof for the exact case

Lemma (4.3). *Let U be a first order critical point of $f(U)$. Then for any $r \times r$ orthonormal matrix R and $\Delta_j = \Delta_j e_j e_j^\top$ ($\Delta = U - U^* R$),*

$$\sum_{j=1}^r \text{vec}(\Delta_j)^\top [\nabla^2 f(U)] \text{vec}(\Delta_j) = \sum_{i=1}^m \left(\sum_{j=1}^r 4 \langle A_i, U \Delta_j^\top \rangle^2 - 2 \langle A_i, U U^\top - U^* U^{*\top} \rangle^2 \right),$$

Proof of Lemma 4.3. For any matrix Z , taking directional second derivative of the function $f(U)$ with respect to Z we get:

$$\begin{aligned} \text{vec}(Z)^\top [\nabla^2 f(U)] \text{vec}(Z) &= \text{vec}(Z)^\top \lim_{t \rightarrow 0} \left[\frac{\nabla f(U + t(Z)) - \nabla f(U)}{t} \right] \\ &= 2 \sum_{i=1}^m \left[2 \langle A_i, U Z^\top \rangle^2 + \langle A_i, U U^\top - U^* U^{*\top} \rangle \langle A_i, Z Z^\top \rangle \right] \end{aligned}$$

Setting $Z = \Delta_j = (U - U^* R) e_j e_j^\top$ and using the first order optimality condition on U , we get,

$$\begin{aligned} &\text{vec}((U - U^* R) e_j e_j^\top)^\top [\nabla^2 f(U)] \text{vec}((U - U^* R) e_j e_j^\top) \\ &= \sum_{i=1}^m 4 \langle A_i, U \Delta_j^\top \rangle^2 + 2 \langle A_i, U U^\top - U^* U^{*\top} \rangle \langle A_i, \Delta_j \Delta_j^\top \rangle \\ &\stackrel{(i)}{=} \sum_{i=1}^m 4 \langle A_i, U e_j e_j^\top \Delta_j^\top \rangle^2 + 2 \langle A_i, U U^\top - U^* U^{*\top} \rangle \langle A_i, U^* e_j e_j^\top (U^* e_j e_j^\top)^\top \rangle \\ &\stackrel{(ii)}{=} \sum_{i=1}^m 4 \langle A_i, U e_j e_j^\top \Delta_j^\top \rangle^2 - 2 \langle A_i, U U^\top - U^* U^{*\top} \rangle \langle A_i, U e_j e_j^\top U^\top - U^* e_j e_j^\top U^{*\top} \rangle. \end{aligned}$$

(i) and (ii) follow from the first order optimality condition (6),

$$\sum_{i=1}^m \langle A_i, U U^\top \rangle U e_j e_j^\top = \sum_{i=1}^m \langle A_i, U^* U^{*\top} \rangle U e_j e_j^\top,$$

for $j = 1 \dots r$. Finally taking sum over j from 1 to r gives the result. \square

Lemma (4.4). *Let U and U^* be two $n \times r$ matrices, and Q is an orthonormal matrix that spans the column space of U . Then there exists an $r \times r$ orthonormal matrix R such that for any first order stationary point U of $f(U)$, the following holds:*

$$\sum_{j=1}^r \|U e_j e_j^\top (U - U^* R)^\top\|_F^2 \leq \frac{1}{8} \|U U^\top - U^* U^{*\top}\|_F^2 + \frac{34}{8} \|(U U^\top - U^* U^{*\top}) Q Q^\top\|_F^2.$$

Proof of Lemma 4.4. To prove this we will expand terms on the both sides in terms of U and $\Delta = U - U^*R$ and then compare. First notice the following properties of R that minimizes $\|U^*R - U\|_F$. Let LSP^\top be the SVD of $U^{*\top}U$. Then, $R = LP^\top$. Hence, $R^\top U^{*\top}U = PSP^\top = U^\top U^*R$ is a PSD matrix. This implies, $U^\top \Delta = U^\top U - U^\top U^*R = U^\top U - R^\top U^{*\top}U = \Delta^\top U$.

Let columns of U be orthogonal, else we can multiply U by an orthonormal matrix and UR will satisfy this. Since UR is also local minimum, and $UU^\top = UR R^\top U^\top$, results for UR will also hold for U . Let Q be the orthonormal matrix that spans the column space of U and $Q_\perp Q_\perp^\top = I - QQ^\top$. Similarly let Q_j span $Ue_j e_j^\top$. Note that Q_j are orthonormal since columns of U are orthogonal. Hence,

$$\begin{aligned} \|(U - U^*R)e_j e_j^\top U^\top\|_F^2 &= \|Ue_j e_j^\top U^\top - Q_j Q_j^\top U^* R e_j e_j^\top U^\top - Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top U^\top\|_F^2 \\ &= \|Ue_j e_j^\top U^\top - Q_j Q_j^\top U^* R e_j e_j^\top U^\top\|_F^2 + \|Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top U^\top\|_F^2 \\ &\leq \frac{\|Ue_j e_j^\top U^\top - Q_j Q_j^\top U^* R e_j e_j^\top (Q_j Q_j^\top U^* R)^\top\|_F^2}{2(\sqrt{2} - 1)} + \|Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top U^\top\|_F^2. \end{aligned} \quad (9)$$

The last inequality follows from Lemma F.1 and the fact that $e_j^\top U^\top U^* R e_j \geq 0, \forall j$ as $U^\top U^* R$ is PSD. Now we will bound the second term in the above equation. The main idea here is to split this term into error between the subspaces of X, X^* and then error between their singular values, since both of them are bounded by distance $\|X - X^* Q Q^\top\|_F$. Let Q^* be an orthonormal matrix that spans the column space of X^* . Also let $X = Q \Sigma_U Q^\top$.

$$\begin{aligned} \|Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top U^\top\|_F^2 &= \text{trace}(e_j^\top R^\top U^{*\top} Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top U^\top U e_j) \\ &= \text{trace}\left(e_j^\top R^\top U^{*\top} Q_{j\perp} Q_{j\perp}^\top U^* R e_j \left[e_j^\top U^\top U e_j - e_j^\top R^\top U^{*\top} Q_j Q_j^\top Q_j Q_j^\top U^* R e_j \right. \right. \\ &\quad \left. \left. + e_j^\top R^\top U^{*\top} Q_j Q_j^\top U^* R e_j \right] \right) \\ &\stackrel{(i)}{\leq} \frac{1}{8} \underbrace{(e_j^\top R^\top U^{*\top} Q_{j\perp} Q_{j\perp}^\top U^* R e_j)^2}_{\text{term1}} + 2 \underbrace{(e_j^\top U^\top U e_j - e_j^\top R^\top U^{*\top} Q_j Q_j^\top Q_j Q_j^\top U^* R e_j)^2}_{\text{term2}} \\ &\quad + \underbrace{(Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top (Q_j Q_j^\top U^* R)^\top)^2}_{\text{term3}}. \end{aligned} \quad (10)$$

where (i) follows from Cauchy-Schwarz inequality.

We will use the following inequality through the rest of the proof. So we state it first for any matrix T .

$$\begin{aligned} \sum_{j=1}^r (e_j^\top T^\top T e_j)^2 &\leq \sum_{j=1}^r \sum_{k=1}^r (e_j^\top T^\top T e_k)^2 \\ &= \sum_{j=1}^r e_j^\top T^\top \left[\sum_{k=1}^r T e_k e_k^\top T^\top \right] T e_j = \sum_{j=1}^r e_j^\top T^\top T T^\top T e_j \\ &= \|T^\top T\|_F^2 = \|T T^\top\|_F^2. \end{aligned} \quad (11)$$

Now we will bound each of the terms in equation .

Term 1: Let, $T = Q_{j\perp} Q_{j\perp}^\top U^* R$. Then applying inequality from equation (11) we get,

$$\begin{aligned} \sum_{j=1}^r (e_j^\top R^\top U^{*\top} Q_{j\perp} Q_{j\perp}^\top U^* R e_j)^2 &= \sum_{j=1}^r (e_j^\top T^\top T e_j)^2 \\ &\leq \|T^\top T\|_F^2 = \|R^\top U^{*\top} Q_\perp Q_\perp^\top U^* R\|_F^2. \end{aligned} \quad (12)$$

Further,

$$\begin{aligned} \|R^\top U^{*\top} Q_\perp Q_\perp^\top U^* R\|_F^2 &= \text{trace}(U^{*\top} Q_\perp Q_\perp^\top U^* U^{*\top} Q_\perp Q_\perp^\top U^*) \\ &= \text{trace}(Q_\perp Q_\perp^\top X^* Q_\perp Q_\perp^\top X^*) \\ &\leq \|Q_\perp Q_\perp^\top X^*\|_F^2 \leq \|X - X^*\|_F^2. \end{aligned} \quad (13)$$

Term 2:

$$\begin{aligned}
& (e_j^\top U^\top U e_j - e_j^\top R^\top U^{*\top} Q_j Q_j^\top U^* R e_j)^2 \\
&= (e_j^\top U^\top U e_j)^2 + (e_j^\top R^\top U^{*\top} Q_j Q_j^\top U^* R e_j)^2 - 2 e_j^\top U^\top U e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top U^* R e_j \\
&= \|U e_j e_j^\top U^\top\|_F^2 + \|Q_j Q_j^\top U^* R e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top\|_F^2 - 2 \text{trace}(e_j^\top U^\top U e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top U^* R e_j) \\
&\stackrel{(i)}{=} \|U e_j e_j^\top U^\top\|_F^2 + \|Q_j Q_j^\top U^* R e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top\|_F^2 - 2 \text{trace}(e_j^\top R^\top U^{*\top} U e_j e_j^\top U^\top U^* R e_j) \\
&= \|U e_j e_j^\top U^\top - Q_j Q_j^\top U^* R e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top\|_F^2. \tag{14}
\end{aligned}$$

(i) follows from $e_j^\top U^\top U e_j = \|U_j\|_F^2$ and $\|U_j\|_F^2 Q_j Q_j^\top = U e_j e_j^\top U^\top$. Now from orthogonality of Q_j we have,

$$\sum_{j=1}^r \|U e_j e_j^\top U^\top - Q_j Q_j^\top U^* R e_j e_j^\top R^\top U^{*\top} Q_j Q_j^\top\|_F^2 \leq \|U U^\top - Q Q^\top U^* U^{*\top} Q Q^\top\|_F^2. \tag{15}$$

Term 3: Finally we bound the last term in equation (10) similar to the first term, which gives,

$$\sum_{j=1}^r (Q_{j\perp} Q_{j\perp}^\top U^* R e_j e_j^\top (Q_j Q_j^\top U^* R)^\top)^2 \leq \|U U^\top - U^* U^{*\top} Q Q^\top\|_F^2.$$

Substituting the above equations (12), (13), (14) and (15) in (9) and (10) gives the result. \square

C Proof for the Noisy Case

In this section we present the proof characterizing the local minima of problem (2). Recall $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{w}$, where \mathbf{X}^* is a rank- r matrix and \mathbf{w} is i.i.d. $\mathcal{N}(0, \sigma_w^2)$.

We consider local optimum that satisfies first and second order optimality conditions of problem (2). In particular \mathbf{U} satisfies $\nabla f(\mathbf{U}) = 0$ and $\mathbf{z}^\top \nabla^2 f(\mathbf{U}) \mathbf{z} \geq 0$ for any $\mathbf{z} \in \mathbb{R}^{n \times r}$. Now we will see how these two conditions constrain the error $\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}$.

C.1 First order optimality

First we will consider the first order condition, $\nabla f(\mathbf{U}) = 0$. For any stationary point \mathbf{U} this implies

$$\sum_i \langle \mathbf{A}_i, \mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \mathbf{A}_i \mathbf{U} = \sum_{i=1}^m w_i \mathbf{A}_i \mathbf{U}. \tag{16}$$

Now using the isometry property of \mathbf{A}_i gives us the following result.

Lemma C.1. [First order condition] For any first order stationary point \mathbf{U} of $f(\mathbf{U})$, and \mathcal{A} satisfying the $(4r, \delta)$ -RIP (3), the following holds:

$$\|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{Q} \mathbf{Q}^\top\|_F \leq \delta \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F + 2 \sqrt{\frac{(1 + \delta) \log(n)}{m}} \sigma_w,$$

w.p. $\geq 1 - \frac{1}{n^2}$, where \mathbf{Q} is an orthonormal matrix that spans the column space of \mathbf{U} .

This lemma states that any stationary point of $f(\mathbf{U})$ is close to a global optimum \mathbf{U}^* in the subspace spanned by columns of \mathbf{U} . Notice that the error along the orthogonal direction $\|\mathbf{X}^* \mathbf{Q}_\perp \mathbf{Q}_\perp^\top\|_F$ can still be large making the distance between \mathbf{X} and \mathbf{X}^* arbitrarily big.

Proof of Lemma C.1. Let $\mathbf{U} = \mathbf{Q} \mathbf{R}$, for some orthonormal \mathbf{Q} . Consider any matrix of the form $\mathbf{Z} \mathbf{Q} \mathbf{R}^\dagger$. The first order optimality condition then implies,

$$\sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{U} \mathbf{R}^\dagger \mathbf{Q}^\top \mathbf{Z}^\top \rangle = \sum_{i=1}^m w_i \mathbf{A}_i \mathbf{U} \mathbf{R}^\dagger \mathbf{Q}^\top \mathbf{Z}^\top.$$

The above equation together with Restricted Isometry Property (equation (5)) gives us the following inequality:

$$\left| \langle \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top \rangle \right| \leq \delta \left\| \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \right\|_F \left\| \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top \right\|_F + 2\sqrt{\frac{(1+\delta)\log(n)}{m}} \sigma_w \left\| \mathbf{Z}^\top \right\|_F,$$

by Cauchy Schwarz inequality and Lemma F.2. Note that for any matrix \mathbf{A} , $\langle \mathbf{A}, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z} \rangle = \langle \mathbf{A}\mathbf{Q}\mathbf{Q}^\top, \mathbf{Z} \rangle$. Furthermore, for any matrix \mathbf{A} , $\sup_{\{\mathbf{Z}: \|\mathbf{Z}\|_F \leq 1\}} \langle \mathbf{A}, \mathbf{Z} \rangle = \|\mathbf{A}\|_F$. Hence the above inequality implies the lemma statement. \square

C.2 Second order optimality

We will now consider the second order condition to show that the error along $\mathbf{Q}_\perp \mathbf{Q}_\perp^\top$ is indeed bounded well. Let $\nabla^2 f(\mathbf{U})$ be the hessian of the objective function. Note that this is an $n \cdot r \times n \cdot r$ matrix. Fortunately for our result we need to only evaluate the Hessian along the direction $\text{vec}(\mathbf{U} - \mathbf{U}^* \mathbf{R})$ for some orthonormal matrix \mathbf{R} .

Lemma C.2. [Hessian computation] Let \mathbf{U} be a first order critical point of $f(\mathbf{U})$. Then for any $r \times r$ orthonormal matrix \mathbf{R} and $\Delta = \mathbf{U} - \mathbf{U}^* \mathbf{R}$,

$$\begin{aligned} & \sum_{j=1}^r \text{vec}(\Delta_j)^\top [\nabla^2 f(\mathbf{U})] \text{vec}(\Delta_j) \\ &= \sum_{i=1}^m \left(\sum_{j=1}^r 4 \langle \mathbf{A}_i, \mathbf{U} \Delta_j^\top \rangle^2 - 2 \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle^2 - 2w_i \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}^* \rangle \right), \end{aligned}$$

Proof of Lemma C.2. For any matrix \mathbf{Z} , taking directional second derivative of the function $f(\mathbf{U})$ with respect to \mathbf{Z} we get:

$$\begin{aligned} & \text{vec}(\mathbf{Z})^\top [\nabla^2 f(\mathbf{U})] \text{vec}(\mathbf{Z}) = \text{vec}(\mathbf{Z})^\top \lim_{t \rightarrow 0} \left[\frac{\nabla f(\mathbf{U} + t(\mathbf{Z})) - \nabla f(\mathbf{U})}{t} \right] \\ &= 2 \sum_{i=1}^m \left[2 \langle \mathbf{A}_i, \mathbf{U} \mathbf{Z}^\top \rangle^2 + \left(\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle - w_i \right) \langle \mathbf{A}_i, \mathbf{Z} \mathbf{Z}^\top \rangle \right] \end{aligned}$$

Setting $\mathbf{Z} = \Delta_j = (\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top$ and using the first order optimality condition on \mathbf{U} , we get,

$$\begin{aligned} & \text{vec}((\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top)^\top [\nabla^2 f(\mathbf{U})] \text{vec}((\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top) \\ &= \sum_{i=1}^m 4 \langle \mathbf{A}_i, \mathbf{U} \Delta_j^\top \rangle^2 + 2 \left(\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle - w_i \right) \langle \mathbf{A}_i, \Delta_j \Delta_j^\top \rangle \\ &= \sum_{i=1}^m 4 \langle \mathbf{A}_i, \mathbf{U} e_j e_j^\top \Delta_j^\top \rangle^2 + 2 \left(\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle - w_i \right) \langle \mathbf{A}_i, \mathbf{U}^* e_j e_j^\top (\mathbf{U}^* e_j e_j^\top)^\top \rangle \\ &= \sum_{i=1}^m 4 \langle \mathbf{A}_i, \mathbf{U} e_j e_j^\top \Delta_j^\top \rangle^2 - 2 \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{U} e_j e_j^\top \mathbf{U}^\top - \mathbf{U}^* e_j e_j^\top \mathbf{U}^{*\top} \rangle \\ &\quad - 2w_i \langle \mathbf{A}_i, \mathbf{U} e_j e_j^\top \mathbf{U}^\top - \mathbf{U}^* e_j e_j^\top \mathbf{U}^{*\top} \rangle. \end{aligned} \tag{17}$$

where the last equality is again by the first order optimality condition (16). \square

Hence from second order optimality of \mathbf{U} we get,

Corollary C.1. [Second order optimality] Let \mathbf{U} be a local minimum of $f(\mathbf{U})$. For any $r \times r$ orthonormal matrix R , w.p. $\geq 1 - \frac{1}{n^2}$,

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^m \left\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \right\rangle^2 &\leq \sum_{i=1}^m \sum_{j=1}^r \left\langle \mathbf{A}_i, \mathbf{U}\Delta_j^\top \right\rangle^2 + \sqrt{\log(n)}\sigma_w \|\mathcal{A}(\mathbf{X} - \mathbf{X}^*)\|_2 \\ &\leq \sum_{i=1}^m \sum_{j=1}^r \left\langle \mathbf{A}_i, \mathbf{U}\Delta_j^\top \right\rangle^2 + 5\log(n)\sigma_w^2 + \frac{1}{20} \sum_{i=1}^m \left\langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}^* \right\rangle^2 \end{aligned}$$

Further for \mathcal{A} satisfying $(2r, \delta)$ -RIP (equation (3)) we have,

$$\frac{1-\delta}{2(1+\delta)} \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F^2 \leq \sum_{j=1}^r \|\mathbf{U}\Delta_j^\top\|_F^2 + \frac{1}{20} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{5\log(n)}{m(1+\delta)} \sigma_w^2. \quad (18)$$

Hence from the above optimality conditions we get the proof of Theorem 4.1.

Proof of Theorem 3.1. Assuming $\mathbf{U}\mathbf{U}^\top \neq \mathbf{U}^*\mathbf{U}^{*\top}$, from Lemma 4.4 and Corollary C.1 we get, with probability $\geq 1 - \frac{2}{n^2}$,

$$\begin{aligned} &\left(\frac{1-\delta}{2(1+\delta)} \right) \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F^2 \\ &\leq \frac{1}{8} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{34}{8} \|\mathbf{X} - \mathbf{X}^* \mathbf{Q} \mathbf{Q}^\top\|_F^2 + \frac{1}{20} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{5\log(n)}{m(1+\delta)} \sigma_w^2 \\ &\stackrel{(i)}{\leq} \left(\frac{1}{8} + \frac{1}{20} \right) \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{34}{8} \left(2\delta^2 \|\mathbf{X} - \mathbf{X}^*\|_F^2 + 8 \frac{(1+\delta)\log(n)}{m} \sigma_w^2 \right) + \frac{5\log(n)}{m(1+\delta)} \sigma_w^2. \end{aligned}$$

(i) follows from Lemma C.1. The above inequality implies,

$$\left(\frac{1-\delta}{2(1+\delta)} - \frac{1}{8} - \frac{1}{20} - \frac{34}{4} \delta^2 \right) \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F^2 \leq 34 \frac{(1+\delta)\log(n)}{m} \sigma_w^2 + \frac{5\log(n)}{m(1+\delta)} \sigma_w^2.$$

If $\delta \leq \frac{1}{10}$, the above inequality reduces to $\|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F \leq c \sqrt{\frac{\log(n)}{m}} \sigma_w$, for some constant $c \leq 17$, w.p. $\geq 1 - \frac{2}{n^2}$. \square

D Proof for the High Rank Case

In this section we will present the proof for the inexact case, where $\text{rank}(\mathbf{X}^*) \geq r$. Recall that measurements are $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$.

Let SVD of \mathbf{X}^* be $\mathbf{Q}^* \Sigma^* \mathbf{Q}^{*\top}$. With slight abuse of notation we use $\mathbf{X}_{jr+1:(j+1)r}^*$ to denote the j th rank r block $\mathbf{Q}_{jr+1:(j+1)r}^{*\top} \Sigma_{jr+1:(j+1)r}^* \mathbf{Q}_{jr+1:(j+1)r}^*$, where $\mathbf{Q}_{jr+1:(j+1)r}^*$ denotes the restriction of \mathbf{Q} to columns $jr+1$ to $(j+1)r$.

D.1 First order optimality

First we will consider the first order condition, $\nabla f(\mathbf{U}) = 0$. For any stationary point \mathbf{U} this implies

$$\sum_i \left\langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \right\rangle \mathbf{A}_i \mathbf{U} = 0. \quad (19)$$

Now using the isometry property of \mathbf{A}_i gives us the following result.

Lemma D.1. [First order condition] For any first order stationary point \mathbf{U} of $f(\mathbf{U})$, and $\{\mathbf{A}_i\}$ satisfying the $(4r, \delta)$ -RIP (3), the following holds:

$$\|\mathbf{X} - \mathbf{Q} \mathbf{Q}^\top \mathbf{X}^*\|_F \leq \delta \|\mathbf{X} - \mathbf{X}_r^*\|_F + \|(\mathbf{X}^* - \mathbf{X}_r^*) \mathbf{Q} \mathbf{Q}^\top\|_F + \delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*$$

where \mathbf{Q} is an orthonormal matrix that spans the column space of \mathbf{U} .

This lemma states that any stationary point of $f(\mathbf{U})$ is close to a global optimum \mathbf{U}^* in the subspace spanned by columns of \mathbf{U} . Notice that the error along the orthogonal direction $\|\mathbf{X}^* \mathbf{Q}_\perp \mathbf{Q}_\perp^\top\|_F$ can still be large making the distance between \mathbf{X} and \mathbf{X}^* arbitrarily big.

Proof of Lemma D.1. Let $\mathbf{U} = \mathbf{Q}\mathbf{R}$, for some orthonormal \mathbf{Q} . Consider any matrix of the form $\mathbf{Z}\mathbf{Q}\mathbf{R}^\top$. The first order optimality condition then implies,

$$\sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle \langle \mathbf{A}_i, \mathbf{U}\mathbf{R}^\top \mathbf{Q}^\top \mathbf{Z}^\top \rangle = \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X}^* - \mathbf{X}_r^* \rangle \langle \mathbf{A}_i, \mathbf{U}\mathbf{R}^\top \mathbf{Q}^\top \mathbf{Z}^\top \rangle.$$

Note that $\mathbf{X} - \mathbf{X}_r^*$ is atmost rank- $2r$. Hence, the above equation together with Restricted Isometry Property (equation (5)) gives us the following inequality:

$$\begin{aligned} |\langle \mathbf{X} - \mathbf{X}_r^*, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top \rangle| - \delta \|\mathbf{X} - \mathbf{X}_r^*\|_F \|\mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top\|_F \\ \leq \frac{1}{m} \sum_{i=1}^m \left\langle \mathbf{A}_i, \sum_j \mathbf{X}_{jr+1:(j+1)r}^* \right\rangle \langle \mathbf{A}_i, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top \rangle \\ \leq \sum_j \left\langle \mathbf{X}_{jr+1:(j+1)r}^*, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z}^\top \right\rangle + \delta \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \\ \leq \|(\mathbf{X}^* - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F + \delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*. \end{aligned}$$

The last inequality follows from $\sum_j \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \leq \|\mathbf{X}^* - \mathbf{X}_r^*\|_*$. The above inequalities are true for any \mathbf{Z} .

Further note that for any matrix \mathbf{A} , $\langle \mathbf{A}, \mathbf{Q}\mathbf{Q}^\top \mathbf{Z} \rangle = \langle \mathbf{A}\mathbf{Q}\mathbf{Q}^\top, \mathbf{Z} \rangle$. Furthermore, for any matrix \mathbf{A} , $\sup_{\{\mathbf{Z}: \|\mathbf{Z}\|_F \leq 1\}} \langle \mathbf{A}, \mathbf{Z} \rangle = \|\mathbf{A}\|_F$. Hence the above inequality implies the Lemma. \square

D.2 Second order optimality

We will now consider the second order condition to show that the error along $\mathbf{Q}_\perp \mathbf{Q}_\perp^\top$ is indeed bounded well. Let $\nabla^2 f(\mathbf{U})$ be the hessian of the objective function. Note that this is an $n \cdot r \times n \cdot r$ matrix. Fortunately for our result we need to only evaluate the Hessian along the direction $\text{vec}(\mathbf{U} - \mathbf{U}^* \mathbf{R})$ for some orthonormal matrix \mathbf{R} .

Lemma D.2. [Hessian computation] Let \mathbf{U} be a first order critical point of $f(\mathbf{U})$. Then for any $n \times r$ matrix \mathbf{Z} ,

$$\text{vec}(\mathbf{Z})^\top [\nabla^2 f(\mathbf{U})] \text{vec}(\mathbf{Z}) = \sum_{i=1}^m 4 \langle \mathbf{A}_i, \mathbf{U}\mathbf{Z}^\top \rangle^2 + 2 \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{Z}\mathbf{Z}^\top \rangle,$$

Further let \mathbf{U} be a local minimum of $f(\mathbf{U})$ and \mathcal{A} satisfying $(2r, \delta)$ -RIP (equation (3)). Then,

$$(1 - 3\delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 \leq 4(1 + \delta) \sum_{j=1}^r \|\mathbf{U} \Delta_j^\top\|_F^2 + \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + \delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2.$$

Proof of Lemma D.2. For any matrix \mathbf{Z} , taking directional second derivative of the function $f(\mathbf{U})$ with respect to \mathbf{Z} we get:

$$\begin{aligned} \text{vec}(\mathbf{Z})^\top [\nabla^2 f(\mathbf{U})] \text{vec}(\mathbf{Z}) &= \text{vec}(\mathbf{Z})^\top \lim_{t \rightarrow 0} \left[\frac{\nabla f(\mathbf{U} + t(\mathbf{Z})) - \nabla f(\mathbf{U})}{t} \right] \\ &= 2 \sum_{i=1}^m \left[2 \langle \mathbf{A}_i, \mathbf{U}\mathbf{Z}^\top \rangle^2 + \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{Z}\mathbf{Z}^\top \rangle \right]. \end{aligned}$$

Setting $\mathbf{Z} = \Delta_j = (\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top$ we get,

$$\begin{aligned}
& \sum_{j=1}^r \text{vec}((\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top)^\top [\nabla^2 f(\mathbf{U})] \text{vec}((\mathbf{U} - \mathbf{U}^* \mathbf{R}) e_j e_j^\top) \\
&= \sum_{i=1}^m \left(\sum_{j=1}^r 4 \langle \mathbf{A}_i, \mathbf{U} e_j e_j^\top (\mathbf{U} - \mathbf{U}_r^* \mathbf{R})^\top \rangle^2 \right. \\
&\quad \left. + 2 \sum_{j=1}^r \langle \mathbf{A}_i, \mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, (\mathbf{U} - \mathbf{U}_r^* \mathbf{R}) e_j e_j^\top (\mathbf{U} - \mathbf{U}_r^* \mathbf{R})^\top \rangle \right) \\
&\stackrel{(i)}{=} \sum_{i=1}^m \left(\sum_{j=1}^r 4 \langle \mathbf{A}_i, \mathbf{U} \Delta_j^\top \rangle^2 + 2 \langle \mathbf{A}_i, \mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{U}_r^* \mathbf{R} (\mathbf{U}_r^* \mathbf{R})^\top - \mathbf{X} \rangle \right).
\end{aligned}$$

(i) is by the first order optimality condition (19).

Hence from second order optimality of \mathbf{U} we get,

$$\sum_{i=1}^m 4 \sum_{j=1}^r \langle \mathbf{A}_i, \mathbf{U} \Delta_j^\top \rangle^2 \geq \sum_{i=1}^m 2 \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}^* \rangle \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle. \quad (20)$$

$$\begin{aligned}
& \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}^* \rangle \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle = \frac{1}{m} \sum_{i=1}^m \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle^2 + \langle \mathbf{A}_i, \mathbf{X}_r^* - \mathbf{X}^* \rangle \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle \\
&\stackrel{(i)}{\geq} (1 - \delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{m} \sum_{i=1}^m \left(\sum_{j=1}^m \langle \mathbf{A}_i, \mathbf{X}_{jr+1:(j+1)r}^* \rangle \right) \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_r^* \rangle \\
&\stackrel{(ii)}{\geq} (1 - \delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \sum_{j=1}^m \langle \mathbf{X} - \mathbf{X}_r^*, \mathbf{X}_{jr+1:(j+1)r}^* \rangle - \delta \sum_{j=1}^m \|\mathbf{X} - \mathbf{X}_r^*\|_F \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \\
&= (1 - \delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \langle \mathbf{X} - \mathbf{X}_r^*, \mathbf{X}^* - \mathbf{X}_r^* \rangle - \delta \sum_{j=1}^m \|\mathbf{X} - \mathbf{X}_r^*\|_F \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \\
&\geq (1 - \delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{2} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 - \delta \sum_{j=1}^m \|\mathbf{X} - \mathbf{X}_r^*\|_F \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \\
&\stackrel{(iii)}{\geq} (1 - \delta) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{2} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 - \delta \frac{1}{2} (\|\mathbf{X} - \mathbf{X}_r^*\|_F^2 + \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2) \\
&= \frac{1 - 3\delta}{2} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 - \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 - \frac{\delta}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2. \quad (21)
\end{aligned}$$

(i) is from using RIP and splitting $\mathbf{X}^* - \mathbf{X}_r^*$ into rank- r components $\mathbf{X}^* - \mathbf{X}_r^* = \sum_{j=1}^{n/r-1} \mathbf{X}_{jr+1:(j+1)r}^*$. (ii) follows from using RIP (5). (iii) follows from $\sum_j \|\mathbf{X}_{jr+1:(j+1)r}^*\|_F \leq \|\mathbf{X}^* - \mathbf{X}_r^*\|_*$.

The Lemma now follows by combining equations (20), (21) and using RIP (3). \square

Hence from the above optimality conditions we get the proof of Theorem 3.4.

Proof of Theorem 3.4. Assuming $\mathbf{U} \mathbf{U}^\top \neq \mathbf{U}_r^* \mathbf{U}_r^{*\top}$, from Lemma 4.4 we know,

$$\sum_{j=1}^r \|\mathbf{U} \Delta_j^\top\|_F^2 \leq \frac{1}{8} \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}_r^* \mathbf{U}_r^{*\top}\|_F^2 + \frac{34}{8} \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}_r^* \mathbf{U}_r^{*\top}) \mathbf{Q} \mathbf{Q}^\top\|_F^2, \quad (22)$$

for some orthonormal \mathbf{R} . Hence combining equations (22), with Lemma D.2 we get,

$$\begin{aligned}
\frac{1 - 3\delta}{2} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 &\leq \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + \frac{\delta}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2 \\
&\quad + 2(1 + \delta) \left(\frac{1}{8} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 + \frac{34}{8} \|(\mathbf{X} - \mathbf{X}_r^*) \mathbf{Q} \mathbf{Q}^\top\|_F^2 \right).
\end{aligned}$$

This implies,

$$\frac{1-7\delta}{4} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 \leq \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + \frac{\delta}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2 + (1+\delta) \frac{17}{2} \|(\mathbf{X} - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F^2. \quad (23)$$

Finally from Lemma D.1 we know,

$$\begin{aligned} \|\mathbf{X} - \mathbf{X}_r^*\mathbf{Q}\mathbf{Q}^\top\|_F^2 &\leq (\delta \|\mathbf{X} - \mathbf{X}_r^*\|_F + \|(\mathbf{X}^* - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F + \delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*)^2 \\ &\leq \frac{11}{10} \|(\mathbf{X}^* - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F^2 + 22\delta^2 \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 + 22\delta^2 \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2. \end{aligned} \quad (24)$$

The last inequality follows from just using $2ab \leq a^2 + b^2$.

Combining equations (23) and (24) gives,

$$\begin{aligned} \left(\frac{1-7\delta}{4} - \frac{17 * 22(1+\delta)\delta^2}{2} \right) \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 &\leq \frac{1}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + \left(\frac{\delta}{2} + \frac{17 * 22\delta^2}{2} \right) \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2 \\ &\quad + (1+\delta) \frac{17 * 11}{20} \|(\mathbf{X}^* - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F^2 \end{aligned}$$

Substituting $\delta = \frac{1}{100}$ gives,

$$\begin{aligned} \|\mathbf{X} - \mathbf{X}_r^*\|_F^2 &\leq \frac{5}{2} \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + 12\delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2 + 10 \|(\mathbf{X}^* - \mathbf{X}_r^*)\mathbf{Q}\mathbf{Q}^\top\|_F^2 \\ &\leq 13 \|\mathbf{X}^* - \mathbf{X}_r^*\|_F^2 + 12\delta \|\mathbf{X}^* - \mathbf{X}_r^*\|_*^2. \end{aligned}$$

□

E Proofs for Section 3

In this section we present the proofs for the strict saddle theorem (Theorem 3.2) and the convergence guarantees (Theorem 3.3). The proofs use the Lemmas developed in Section 4 and the supporting Lemmas from Section F.

Proof of Theorem 3.2. From Lemma 4.3 we know that

$$\begin{aligned} \sum_{j=1}^r \text{vec}(\Delta_j)^\top \left[\frac{1}{m} \nabla^2 f(\mathbf{U}) \right] \text{vec}(\Delta_j) &= \frac{1}{m} \sum_{i=1}^m \left(\sum_{j=1}^r 4 \langle \mathbf{A}_i, \mathbf{U} \Delta_j^\top \rangle^2 - 2 \langle \mathbf{A}_i, \mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top} \rangle \right)^2 \\ &\leq 4(1+\delta) \sum_{j=1}^r \|\mathbf{U} \Delta_j^\top\|_F^2 - 2(1-\delta) \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2, \end{aligned} \quad (25)$$

where the last inequality follows from the RIP (3). Now applying Lemma 4.4 in equation (25) we get,

$$\begin{aligned} \sum_{j=1}^r \text{vec}(\Delta_j)^\top \left[\frac{1}{m} \nabla^2 f(\mathbf{U}) \right] \text{vec}(\Delta_j) &\leq (1+\delta) \left(\frac{1}{2} \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 + 17 \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top})\mathbf{Q}\mathbf{Q}^\top\|_F^2 \right) - 2(1-\delta) \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 \\ &= 17(1+\delta) \|(\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top})\mathbf{Q}\mathbf{Q}^\top\|_F^2 - \frac{(3-5\delta)}{2} \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 \\ &\stackrel{(i)}{\leq} \left[17(1+\delta)\delta^2 - \frac{(3-5\delta)}{2} \right] \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2 \\ &\stackrel{(ii)}{\leq} -1 \cdot \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2. \end{aligned} \quad (26)$$

(i) follows from Lemma 4.2. (ii) follows from $\delta \leq 1/10$. Now notice that from lemma F.1

$$\begin{aligned}\|\mathbf{X} - \mathbf{X}^*\|_F^2 &\geq 2(\sqrt{2} - 1)\|(\mathbf{U} - \mathbf{U}^*R)(\mathbf{U}^*R)^\top\|_F^2 \\ &\geq 2(\sqrt{2} - 1)\sigma_r(\mathbf{X}^*)\|\mathbf{U} - \mathbf{U}^*R\|_F^2.\end{aligned}\quad (27)$$

Finally notice that $\Delta_j = \Delta e_j e_j^\top$ have orthogonal columns. Hence,

$$\begin{aligned}\lambda_{\min}\left[\frac{1}{m}\nabla^2(f(\mathbf{U}, \mathbf{V}))\right] &\leq \frac{1}{\|\mathbf{U} - \mathbf{U}^*R\|_F^2} \sum_{j=1}^r \text{vec}(\Delta_j)^\top \left[\frac{1}{m}\nabla^2 f(\mathbf{U})\right] \text{vec}(\Delta_j) \\ &\stackrel{(i)}{\leq} \frac{-1}{\|\mathbf{U} - \mathbf{U}^*R\|_F^2} \left\|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\right\|_F^2 \\ &\stackrel{(ii)}{\leq} \frac{-2(\sqrt{2} - 1)\sigma_r(\mathbf{X}^*)\|\mathbf{U} - \mathbf{U}^*R\|_F^2}{\|\mathbf{U} - \mathbf{U}^*R\|_F^2} \\ &\leq \frac{-4}{5}\sigma_r(\mathbf{X}^*).\end{aligned}$$

(i) follows from equation (26). (ii) follows from equation (27). \square

Proof of Theorem 3.3. To prove this theorem we use Theorem 6 of Ge et al. [10]. We need to show that $f(\mathbf{U})$ satisfies, 1) strict saddle property, 2) local strong convexity, 3) f is bounded, smooth and has Lipschitz Hessian.

The boundedness assumption easily follows from assuming we are optimizing over a bounded domain b such that, $\|\mathbf{U}^*\|_F \leq b$. Note that we can have any reasonable upper bound on the optimum and we can easily estimate this from $\sum_i y_i^2$ which is $\geq (1 - \delta)\|\mathbf{X}^*\|_F^2$ for the noiseless case.

Finally all the calculations below are for scaled version of $f(x)$ by $\frac{1}{m}$. Note that this does not change the number of iterations as both smoothness and strong convexity parameters are scaled by the same constant.

Smoothness constant β : Recall that smoothness of f is bounded by maximum eigenvalue of Hessian over the domain. Hence, $\beta = \max_{\mathbf{Z}: \|\mathbf{Z}\|_F \leq 1} \mathbf{Z}^\top \nabla^2 f(\mathbf{U}) \mathbf{Z}$. We have computed this projection of Hessian in Lemma C.2. Hence,

$$\begin{aligned}\beta &= 2 \max_{\mathbf{Z}: \|\mathbf{Z}\|_F^2 \leq 1} \sum_{i=1}^m \left[2\langle \mathbf{A}_i, \mathbf{U}\mathbf{Z}^\top \rangle^2 + \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{Z}\mathbf{Z}^\top \rangle \right] \\ &\stackrel{(i)}{\leq} \max_{\mathbf{Z}: \|\mathbf{Z}\|_F^2 \leq 1} 2(2(1 + \delta)\|\mathbf{U}\|_F^2\|\mathbf{Z}\|_F^2 + (1 + \delta)\|\mathbf{X} - \mathbf{X}^*\|_F\|\mathbf{Z}\mathbf{Z}^\top\|_F) \\ &\leq 4(1 + \delta)b^2 + (1 + \delta)2b \leq 5b^2 + 3b.\end{aligned}$$

(i) follows from the RIP.

ρ -Lipschitz Hessian: Now we will compute the Lipschitz constant of Hessian of $f(\mathbf{U})$. We will first bound the spectral norm of difference of Hessian at two points \mathbf{U}, \mathbf{V} in terms of $\|\mathbf{U} - \mathbf{V}\|_F$ along orthogonal direction \mathbf{Z}_i and combine them to get bound on ρ . Given two $n \times r$ matrices \mathbf{U}, \mathbf{V} ,

$$\begin{aligned}&\langle \nabla^2 f(\mathbf{U}) - \nabla^2 f(\mathbf{V}), \mathbf{Z}\mathbf{Z}^\top \rangle \\ &\leq 2 \max_{\mathbf{Z}: \|\mathbf{Z}\|_F^2 \leq 1} \sum_{i=1}^m \left[2\langle \mathbf{A}_i, \mathbf{U}\mathbf{Z}^\top \rangle^2 + \langle \mathbf{A}_i, \mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{Z}\mathbf{Z}^\top \rangle \right] \\ &\quad - \sum_{i=1}^m \left[2\langle \mathbf{A}_i, \mathbf{V}\mathbf{Z}^\top \rangle^2 + \langle \mathbf{A}_i, \mathbf{V}\mathbf{V}^\top - \mathbf{U}^*\mathbf{U}^{*\top} \rangle \langle \mathbf{A}_i, \mathbf{Z}\mathbf{Z}^\top \rangle \right] \\ &\leq 4(1 + \delta)(\|\mathbf{U}\mathbf{Z}^\top\|_F^2 - \|\mathbf{V}\mathbf{Z}^\top\|_F^2) + 2(1 + \delta)\|\mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top\|_F\|\mathbf{Z}\mathbf{Z}^\top\|_F \\ &\leq 4(1 + \delta)\|\mathbf{Z}\|_F^2(\|\mathbf{U} - \mathbf{V}\|_F^2 + 2\|\mathbf{U}\|_F\|\mathbf{U} - \mathbf{V}\|_F) + 2(1 + \delta)\|\mathbf{U}\mathbf{U}^\top - \mathbf{V}\mathbf{V}^\top\|_F \\ &\leq \|\mathbf{Z}\|_F^2\|\mathbf{U} - \mathbf{V}\|_F(8(1 + \delta)b + 4(1 + \delta)b) \\ &= \|\mathbf{Z}\|_F^2\|\mathbf{U} - \mathbf{V}\|_F(12(1 + \delta)b).\end{aligned}\quad (28)$$

Hence, using the variational characterization of the Frobenius norm, the Hessian Lipschitz constant is bounded by $\max \{ \mathbf{Z}_i \} \sum_i \langle \nabla^2 f(\mathbf{U}) - \nabla^2 f(\mathbf{V}), \mathbf{Z}_i \mathbf{Z}_i^\top \rangle$, where \mathbf{Z}_i are orthogonal with $\sum_i \|\mathbf{Z}_i\|_F^2 \leq 1$. Hence from equation (28) we get $\rho = O(b)$.

Strict saddle property: So far we have shown regularity properties of $f(\mathbf{U})$. Now we will discuss the strict saddle property. Theorem 3.2 shows that $\lambda_{\min} [\nabla^2(f(\mathbf{U}))] \leq \frac{-2}{5} \sigma_r(\mathbf{X}^*)$. To use results of [10] we need to show this property over an ϵ neighborhood of any saddle point \mathbf{U} . For this first recall by smoothness, $\|\nabla f(\mathbf{U}) - \nabla f(\mathbf{V})\|_F \leq \beta \|\mathbf{U} - \mathbf{V}\|_F$. Therefore $\nabla f(\mathbf{V}) \leq \epsilon$, when $\|\mathbf{U} - \mathbf{V}\|_F \leq \frac{\epsilon}{\beta}$. Further we know the Hessian spectral norm is ρ Lipschitz from equation (28). Hence, for any direction \mathbf{Z} ,

$$\mathbf{Z}^\top (\nabla^2(f(\mathbf{V})) - \nabla^2(f(\mathbf{U}))) \mathbf{Z}^\top \leq \rho \|\mathbf{U} - \mathbf{V}\|_F \leq \rho \frac{\epsilon}{\beta}.$$

In particular choosing \mathbf{Z} to be the projection direction, $\mathbf{U} - \mathbf{U}^*$ implies from Theorem 3.2,

$$\mathbf{Z}^\top (\nabla^2(f(\mathbf{V}))) \mathbf{Z}^\top \leq \frac{-2}{5} \sigma_r(\mathbf{X}^*) + \rho \frac{\epsilon}{\beta}.$$

Hence for all \mathbf{V} in the bowl of radius ϵ around \mathbf{U} , where $\epsilon \leq \frac{\beta}{5\rho} \sigma_r(\mathbf{X}^*)$,

$$\lambda_{\min} [\nabla^2(f(\mathbf{V}))] \leq \frac{-1}{5} \sigma_r(\mathbf{X}^*). \quad (29)$$

Local strong convexity: Finally we need to show that the function is α strongly convex in a neighborhood θ around the optimum $\mathbf{U}^* R$, for any orthonormal R . This easily follows from existing local convergence results for this problem. For example, Lemma 6.1 of Bhojanapalli et al. [2] states that, for $\|\mathbf{U} - \mathbf{U}^* R\|_F \leq \frac{\sigma_r(\mathbf{X}^*)}{200\sigma_1(\mathbf{X}^*)} \sigma_r(\mathbf{U}^* R)$,

$$\langle \nabla f(\mathbf{U}), \mathbf{U} - \mathbf{U}^* R \rangle \geq \frac{2}{3} \eta \|\nabla f(\mathbf{U})\|_F^2 + \frac{27}{200} \sigma_r(\mathbf{U}^* R)^2 \|\mathbf{U} - \mathbf{U}^* R\|_F^2. \quad (30)$$

for $\delta = \frac{1}{10}$ and some step size $\eta \propto \frac{1}{\|\mathbf{X}^*\|_2}$. Hence $f(\mathbf{U})$ is locally strong convex with $\alpha = \frac{27}{200} \sigma_r(\mathbf{U}^* R)^2$ in the neighborhood of radius $\theta = \frac{\sigma_r(\mathbf{X}^*)}{200\sigma_1(\mathbf{X}^*)} \sigma_r(\mathbf{U}^* R)$ around the optimum.

Substituting these parameters in the Theorem 6 of Ge et al. [10] gives the result. \square

F Supporting Lemmas

In this section we present the supporting results used in the proofs above.

The following lemma relates the error $\|(\mathbf{U} - \mathbf{Y})\mathbf{U}^\top\|_F$ with $\|\mathbf{U}\mathbf{U}^\top - \mathbf{Y}\mathbf{Y}^\top\|_F$ under some conditions on \mathbf{U} and \mathbf{Y} . This is a generalization of Lemma 5.4 in [26] and the proof follows similarly.

Lemma F.1. *Let \mathbf{U} and \mathbf{Y} be two $n \times r$ matrices. Further let $\mathbf{U}^\top \mathbf{Y} = \mathbf{Y}^\top \mathbf{U}$ be a PSD matrix. Then,*

$$\|(\mathbf{U} - \mathbf{Y})\mathbf{U}^\top\|_F^2 \leq \frac{1}{2(\sqrt{2} - 1)} \|\mathbf{U}\mathbf{U}^\top - \mathbf{Y}\mathbf{Y}^\top\|_F^2.$$

Proof. To prove this we will expand terms on the both sides in terms of \mathbf{U} and $\Delta = \mathbf{U} - \mathbf{Y}$ and then compare.

$$\begin{aligned}
& \|(\mathbf{U}\mathbf{U}^\top - \mathbf{Y}\mathbf{Y}^\top)\|_F^2 = \|(\mathbf{U}\Delta^\top + \Delta\mathbf{U}^\top - \Delta\Delta^\top)\|_F^2 \\
& = \text{trace}(\Delta\mathbf{U}^\top\mathbf{U}\Delta^\top + \mathbf{U}\Delta^\top\Delta\mathbf{U}^\top + \Delta\Delta^\top\Delta\Delta^\top + 2\Delta\mathbf{U}^\top\Delta\mathbf{U}^\top - 2\Delta\Delta^\top\Delta\mathbf{U}^\top - 2\Delta\Delta^\top\mathbf{U}\Delta^\top) \\
& \stackrel{(i)}{=} \text{trace}(2\mathbf{U}^\top\mathbf{U}\Delta^\top\Delta + (\Delta^\top\Delta)^2 + 2(\mathbf{U}^\top\Delta)^2 - 4\Delta^\top\Delta\mathbf{U}^\top\Delta) \\
& \stackrel{(ii)}{=} \text{trace}(2\mathbf{U}^\top\mathbf{U}\Delta^\top\Delta + (\Delta^\top\Delta - \sqrt{2}\mathbf{U}^\top\Delta)^2 - 2(2 - \sqrt{2})\Delta^\top\Delta\mathbf{U}^\top\Delta) \\
& \stackrel{(iii)}{\geq} 2 \text{trace}\left(\left[\mathbf{U}^\top\mathbf{U} - (2 - \sqrt{2})\mathbf{U}^\top\Delta\right]\Delta^\top\Delta\right) \\
& = 2 \text{trace}\left(\left[(\sqrt{2} - 1)\mathbf{U}^\top\mathbf{U} + (2 - \sqrt{2})\mathbf{U}^\top\mathbf{Y}\right]\Delta^\top\Delta\right) \\
& \stackrel{(iv)}{\geq} 2 \text{trace}\left((\sqrt{2} - 1)\mathbf{U}^\top\mathbf{U}\Delta^\top\Delta\right).
\end{aligned}$$

(i) follows from the following properties of trace: $\text{trace}(\mathbf{A}\mathbf{B}) = \text{trace}(\mathbf{B}\mathbf{A})$ and $\text{trace}(\mathbf{A}) = \text{trace}(\mathbf{A}^\top)$. (ii) follows from completing the squares. (iii) follows from $\text{trace}(\mathbf{A}^2) \geq 0$. (iv) follows from the hypothesis of the lemma ($\mathbf{U}^\top\mathbf{Y}$ is PSD) and $\text{trace}(\mathbf{A}\mathbf{B}) \geq 0$ for PSD matrices \mathbf{A} and \mathbf{B} .

Finally notice that $\|(\mathbf{U} - \mathbf{Y})\mathbf{U}^\top\|_F^2 = \text{trace}(\mathbf{U}^\top\mathbf{U}\Delta^\top\Delta)$. This completes the proof. \square

We recall the standard Gaussian random variable concentration here.

Lemma F.2. *Let $w_i \approx \mathcal{N}(0, \sigma_w)$, then*

$$\sum_{i=1}^m w_i x_i \leq 2\sqrt{\log(n)}\sigma_w\|\mathbf{x}\|,$$

with probability $\geq 1 - \frac{1}{n^2}$.

Proof. Recall $\mathbb{E}[e^{tw_i}] = e^{\sigma_w^2 t^2/2}$. Then by Markov's inequality, $P(\sum_{i=1}^m w_i x_i \geq c\|\mathbf{x}\|) \leq \frac{e^{\sigma_w^2 \|\mathbf{x}\|^2 t^2/2}}{e^{tc\|\mathbf{x}\|}} \leq e^{-c^2/2\sigma_w^2}$, by setting $t = \frac{c}{\sigma_w^2 \|\mathbf{x}\|}$. Choosing $c = 2\sqrt{\log(n)}\sigma_w$ completes the proof. \square