

# Supplementary Material

## Martingale

A filtration, defined on a measurable probability space, is an increasing sequence of sub-sigma algebras  $\{\mathcal{F}_t\}$  for  $t \geq 0$ , meaning that  $\mathcal{F}_s \subset \mathcal{F}_t$  for all  $s \leq t$ . In our context,  $\mathcal{F}_t$  encodes the set of all the random variables seen thus far (i.e., from 0 to  $t$ ). For brevity,  $\mathbb{E}[\cdot | \mathbf{X}_t] \triangleq \mathbb{E}[\cdot | \mathcal{F}_t]$ . Let  $H = \{H_t\}$  and  $\mathcal{F} = \{\mathcal{F}_t\}$  be a stochastic process and a filtration on the same probability space, respectively.  $H$  then is called a martingale (super-martingale) with respect to  $\mathcal{F}$  if for each  $t$ ,  $H_t$  is  $\mathcal{F}_t$ -measurable,  $\mathbb{E}[|H_t|] < \infty$ , and  $\mathbb{E}[H_{t+1} | \mathcal{F}_t] = H_t$  ( $\mathbb{E}[H_{t+1} | \mathcal{F}_t] \leq H_t$ ). Given a random variable  $X \geq 0$  and a constant  $a > 0$ , we have the Markov inequality  $P(X \geq a) \leq \mathbb{E}[X]/a$ . Let  $X_0, X_1, \dots, X_T$  be a martingale or supermartingale such that  $|X_t - X_{t-1}| \leq d_t$  (i.e., bounded difference) where  $d_t$  is a deterministic function of  $t$ . We then have the Azuma-Hoeffding (Mitzenmacher and Upfal, 2005) inequality  $P(X_t - X_0 \geq a) \leq \exp\{-a^2/(2 \sum_{s=1}^t d_s^2)\}$  for all  $t \geq 0$  and any  $a > 0$ .

## von Neumann's Trace Inequality

For any two real and symmetric  $n \times n$  matrices  $\mathbf{B}$  and  $\mathbf{C}$  with eigenvalues indexed in descending order, it holds that

$$\max \left\{ \sum_{i=1}^n \lambda_{n-i+1}(\mathbf{B}) \lambda_1(\mathbf{C}), \sum_{i=1}^n \lambda_i(\mathbf{B}) \lambda_{n-i+1}(\mathbf{C}) \right\} \leq \text{tr}(\mathbf{B}\mathbf{C}) \leq \sum_{i=1}^n \lambda_i(\mathbf{B}) \lambda_i(\mathbf{C}).$$

*Proof.* See (Lewis, 1996) for the right inequality. The proof of the left one is completed by replacing  $\mathbf{B}$  with  $-\mathbf{B}$  or replacing  $\mathbf{C}$  with  $-\mathbf{C}$  in the right one.  $\square$

**Lemma 5.1.** Let  $\tilde{\beta} = \max_i \|\tilde{\mathbf{A}}_i - \mathbf{A}\|_2$ . Then

$$\begin{aligned} \|\mathbf{W}_t\|_2 &\leq \mu \triangleq 4\tilde{\beta}, \\ \|\mathbf{W}_t\|_F^2 &\leq \nu_t^2 \triangleq 24\tilde{\beta}^2 \left( \Theta(\mathbf{X}_t, \mathbf{U}) + \Theta(\tilde{\mathbf{X}}_{s-1}, \mathbf{U}) \right). \end{aligned}$$

*Proof.* For brevity, let's omit subscripts with  $\mathbf{X}$ ,  $\tilde{\mathbf{X}}$ ,  $\mathbf{W}$  and  $\mathbf{Q}$  herein.  $\mathbf{W}$  then can be written as

$$\begin{aligned} \mathbf{W} &= \mathbf{X}_\perp \mathbf{X}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) (\mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q}) + \mathbf{X}_\perp \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \\ &\quad - \mathbf{X} \text{skew} \left( \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \tilde{\mathbf{X}}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right). \end{aligned}$$

We then have

$$\begin{aligned} \|\mathbf{W}\|_2 &\leq \left\| \mathbf{X}_\perp \mathbf{X}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) (\mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q}) \right\|_2 + \left\| \mathbf{X}_\perp \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right\|_2 \\ &\quad + \left\| \mathbf{X} \text{skew} \left( \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \tilde{\mathbf{X}}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right) \right\|_2 \\ &\leq \tilde{\beta} \left( \|\mathbf{X}\|_2 + \|\tilde{\mathbf{X}}\mathbf{Q}\|_2 \right) + 2\tilde{\beta} = 4\tilde{\beta}, \end{aligned}$$

and

$$\begin{aligned}
\|\mathbf{W}\|_F^2 &\leq \left( \left\| \mathbf{X}_\perp \mathbf{X}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) (\mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q}) \right\|_F + \left\| \mathbf{X}_\perp \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \tilde{\mathbf{X}}^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right\|_F \right. \\
&\quad \left. + \left\| \mathbf{X}_{\text{skew}} \left( \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \tilde{\mathbf{X}}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right) \right\|_F \right)^2 \\
&\leq \left( \left\| \mathbf{X}_\perp \mathbf{X}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) \right\|_2 \left\| \mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q} \right\|_F + \left\| \mathbf{X}_\perp \right\|_2 \left\| \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \right\|_F \left\| \tilde{\mathbf{X}}^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right\|_2 \right. \\
&\quad \left. + \left\| \mathbf{X} \right\|_2 \left\| \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \right\|_F \left\| \tilde{\mathbf{X}}_\perp^\top (\mathbf{A}_{t+1} - \mathbf{A}) \tilde{\mathbf{X}}\mathbf{Q} \right\|_2 \right)^2 \\
&\leq \tilde{\beta}^2 \left( \left\| \mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q} \right\|_F + \left\| \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \right\|_F + \left\| \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \right\|_F \right)^2 \\
&\leq 3\tilde{\beta}^2 \left( \left\| \mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q} \right\|_F^2 + \left\| \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \right\|_F^2 + \left\| \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \right\|_F^2 \right).
\end{aligned}$$

From the proof of Lemma 11 in (Shamir, 2016), we have  $\left\| \mathbf{X} - \tilde{\mathbf{X}}\mathbf{Q} \right\|_F^2 \leq 4 \left( \Theta(\mathbf{X}, \mathbf{U}) + \Theta(\tilde{\mathbf{X}}, \mathbf{U}) \right)$ . Further noting that  $\mathbf{I} = \mathbf{U}\mathbf{U}^\top + \mathbf{U}_\perp \mathbf{U}_\perp^\top = \mathbf{X}\mathbf{X}^\top + \mathbf{X}_\perp \mathbf{X}_\perp^\top$ , we have

$$\begin{aligned}
\left\| \mathbf{X}_\perp^\top \tilde{\mathbf{X}} \right\|_F^2 &= \left\| \mathbf{X}_\perp^\top (\mathbf{U}\mathbf{U}^\top + \mathbf{U}_\perp \mathbf{U}_\perp^\top) \tilde{\mathbf{X}} \right\|_F^2 \leq \left( \left\| \mathbf{X}_\perp^\top \mathbf{U}\mathbf{U}^\top \tilde{\mathbf{X}} \right\|_F + \left\| \mathbf{X}_\perp^\top \mathbf{U}_\perp \mathbf{U}_\perp^\top \tilde{\mathbf{X}} \right\|_F \right)^2 \\
&\leq \left( \left\| \mathbf{X}_\perp^\top \mathbf{U} \right\|_F \left\| \mathbf{U}^\top \tilde{\mathbf{X}} \right\|_2 + \left\| \mathbf{X}_\perp^\top \mathbf{U}_\perp \right\|_2 \left\| \mathbf{U}_\perp^\top \tilde{\mathbf{X}} \right\|_F \right)^2 \leq \left( \left\| \mathbf{X}_\perp^\top \mathbf{U} \right\|_F + \left\| \mathbf{U}_\perp^\top \tilde{\mathbf{X}} \right\|_F \right)^2 \\
&\leq 2 \left( \left\| \mathbf{X}_\perp^\top \mathbf{U} \right\|_F^2 + \left\| \mathbf{U}_\perp^\top \tilde{\mathbf{X}} \right\|_F^2 \right) \\
&= 2 \left( k - \left\| \mathbf{U}^\top \mathbf{X} \right\|_F^2 + k - \left\| \mathbf{U}^\top \tilde{\mathbf{X}} \right\|_F^2 \right) = 2 \left( \Theta(\mathbf{X}, \mathbf{U}) + \Theta(\tilde{\mathbf{X}}, \mathbf{U}) \right),
\end{aligned}$$

and similarly  $\left\| \mathbf{X}^\top \tilde{\mathbf{X}}_\perp \right\|_F^2 \leq 2 \left( \Theta(\mathbf{X}, \mathbf{U}) + \Theta(\tilde{\mathbf{X}}, \mathbf{U}) \right)$ . Therefore, we arrive at

$$\|\mathbf{W}\|_F^2 \leq 24\tilde{\beta}^2 \left( \Theta(\mathbf{X}, \mathbf{U}) + \Theta(\tilde{\mathbf{X}}, \mathbf{U}) \right).$$

□

**Lemma 5.2.** If  $\mathbf{C}_1 \succ \mathbf{0}$ ,  $\mathbf{C}_2 \succ \mathbf{0}$  and  $\mathbf{D}_2 \succ \mathbf{D}_1 \succ \mathbf{0}$ , then  $\text{tr}(\mathbf{C}_1 \mathbf{D}_1^{-1}) \geq \text{tr}((\mathbf{C}_1 - \mathbf{C}_2) \mathbf{D}_2^{-1})$ .

*Proof.*

$$\begin{aligned}
\mathbf{C}_1 \mathbf{D}_1^{-1} &= \mathbf{C}_1 (\mathbf{D}_2 - (\mathbf{D}_2 - \mathbf{D}_1))^{-1} \\
&= \mathbf{C}_1 \mathbf{D}_2^{-1/2} \left( \mathbf{I} - \mathbf{D}_2^{-1/2} (\mathbf{D}_2 - \mathbf{D}_1) \mathbf{D}_2^{-1/2} \right)^{-1} \mathbf{D}_2^{-1/2}.
\end{aligned}$$

By Lemmas 2 and 3 in (Shamir, 2016), we have

$$\begin{aligned}
\text{tr}(\mathbf{C}_1 \mathbf{D}_1^{-1}) &= \text{tr} \left( \mathbf{C}_1 \mathbf{D}_2^{-1/2} (\mathbf{I} - \mathbf{D}_2^{-1/2} (\mathbf{D}_2 - \mathbf{D}_1) \mathbf{D}_2^{-1/2})^{-1} \mathbf{D}_2^{-1/2} \right) \\
&= \text{tr} \left( \mathbf{D}_2^{-1/2} \mathbf{C}_1 \mathbf{D}_2^{-1/2} \left( \mathbf{I} - \mathbf{D}_2^{-1/2} (\mathbf{D}_2 - \mathbf{D}_1) \mathbf{D}_2^{-1/2} \right)^{-1} \right) \\
&\geq \text{tr} \left( \mathbf{D}_2^{-1/2} \mathbf{C}_1 \mathbf{D}_2^{-1/2} \left( \mathbf{I} + \mathbf{D}_2^{-1/2} (\mathbf{D}_2 - \mathbf{D}_1) \mathbf{D}_2^{-1/2} \right) \right) \\
&\geq \text{tr} \left( \mathbf{D}_2^{-1/2} \mathbf{C}_1 \mathbf{D}_2^{-1/2} \right) = \text{tr}(\mathbf{C}_1 \mathbf{D}_2^{-1}) \geq \text{tr}((\mathbf{C}_1 - \mathbf{C}_2) \mathbf{D}_2^{-1}).
\end{aligned}$$

□

**Lemma 5.3.** For  $i = 1, 2$  and any  $\varsigma \in [0, 1]$ ,

$$\begin{aligned} \|b_i(\mathbf{W}_t)\|_F &\leq 2(1 + \alpha\beta)\alpha\nu_t, \\ \|a_1(\mathbf{X}_t)\|_2 + \|b_1(\mathbf{W}_t)\|_2 &\leq (1 + \alpha\beta)(1 + \alpha(\beta + 2\mu)), \\ a_2(\mathbf{X}_t) + \varsigma b_2(\mathbf{W}_t) &\succcurlyeq (1 - 2(1 + \alpha\beta)\alpha\mu) \mathbf{I}. \end{aligned}$$

*Proof.* We omit subscripts for brevity herein. Since  $\|\mathbf{Y}\|_2 = \|(\mathbf{I} + \alpha\mathbf{X}_\perp\mathbf{X}_\perp^\top\mathbf{A})\mathbf{X}\|_2 \leq 1 + \alpha\beta$ , we have

$$\begin{aligned} \|b_1(\mathbf{W})\|_F &= 2\alpha \|\text{sym}(\mathbf{Y}^\top \mathbf{U} \mathbf{U}^\top \mathbf{W})\|_F \leq 2\alpha \|\mathbf{Y}^\top \mathbf{U} \mathbf{U}^\top \mathbf{W}\|_F \leq 2\alpha(1 + \alpha\beta)\nu_t, \\ \|b_2(\mathbf{W})\|_F &= 2\alpha \|\text{sym}(\mathbf{Y}^\top \mathbf{W})\|_F \leq 2\alpha \|\mathbf{Y}^\top \mathbf{W}\|_F \leq 2\alpha(1 + \alpha\beta)\nu_t. \end{aligned}$$

In addition, we can write

$$\|a_1(\mathbf{X})\|_2 + \|b_1(\mathbf{W})\|_2 \leq \|\mathbf{Y}^\top \mathbf{U}\|_2 (\|\mathbf{Y}^\top \mathbf{U}\|_2 + 2\alpha \|\mathbf{U}^\top \mathbf{W}\|_2) \leq (1 + \alpha\beta)(1 + \alpha(\beta + 2\mu)).$$

Further we note that  $\mathbf{Y}^\top \mathbf{Y} = \mathbf{I} + \alpha\mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp^\top \mathbf{X}_\perp \mathbf{A} \mathbf{X}$ .

$$a_2(\mathbf{X}) + \varsigma b_1(\mathbf{W}) \succcurlyeq (1 - \|b_1(\mathbf{W})\|_2) \mathbf{I} \succcurlyeq (1 - 2\alpha(1 + \alpha\beta)\mu) \mathbf{I}.$$

□

**Lemma 5.4.**  $\text{tr}(\mathbf{X}_t^\top \mathbf{U} \mathbf{U}^\top (\mathbf{I} - \mathbf{X}_t \mathbf{X}_t^\top) \mathbf{A} \mathbf{X}) \geq \tau \left( \|\mathbf{X}_t^\top \mathbf{U}\|_F^2 - \|\mathbf{X}_t^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_t\|_F^2 \right).$

*Proof.* Since  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top + \mathbf{U}_\perp \mathbf{\Sigma}_\perp \mathbf{U}_\perp^\top$ , then

$$\begin{aligned} &\text{tr}(\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) \\ &= \text{tr}(\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \mathbf{\Sigma} \mathbf{U}^\top \mathbf{X}) + \text{tr}(\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \mathbf{\Sigma}_\perp \mathbf{U}_\perp^\top \mathbf{X}) \\ &= \text{tr}(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \mathbf{\Sigma}) + \text{tr}(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \mathbf{\Sigma}_\perp). \end{aligned}$$

By von Neumann's trace inequality, we have

$$\begin{aligned} \text{tr}(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \mathbf{\Sigma}) &\geq \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) \lambda_{k-i+1}(\mathbf{\Sigma}), \\ \text{tr}(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \mathbf{\Sigma}_\perp) &\geq \sum_{i=1}^{n-k} \lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) \lambda_{n-k-i+1}(\mathbf{\Sigma}_\perp). \end{aligned}$$

Note that both matrices above, i.e.,  $\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}$  and  $\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp$ , are symmetric and thus von Neumann's trace inequality can be applied. In fact,

$$\begin{aligned} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} &= \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top (\mathbf{I} - \mathbf{X} \mathbf{X}^\top) \mathbf{U} \\ &= \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} - \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \\ &= \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} - (\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U})^2, \end{aligned}$$

which is symmetric. Furthermore, it is positive semi-definite, because

$$\|\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}\|_2 \leq (\|\mathbf{X}^\top\|_2 \|\mathbf{U}\|_2)^2 = 1$$

and thus

$$\lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) = \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}) - \lambda_i^2(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}) \geq 0.$$

Likewise,

$$\begin{aligned} \mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp &= \mathbf{U}_\perp^\top (\mathbf{I} - \mathbf{X}_\perp \mathbf{X}_\perp^\top) \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \\ &= -\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top (\mathbf{I} - \mathbf{U}_\perp \mathbf{U}_\perp^\top) \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \\ &= -\left( \mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp - (\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp)^2 \right), \end{aligned}$$

which is symmetric but negative semi-definite now, because

$$\|\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp\|_2 \leq (\|\mathbf{X}_\perp^\top\|_2 \|\mathbf{U}_\perp\|_2)^2 = 1$$

and thus

$$\lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) = -(\lambda_i(\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) - \lambda_i^2(\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp)) \leq 0.$$

We now can write

$$\begin{aligned} & \text{tr}(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \boldsymbol{\Sigma}) + \text{tr}(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \boldsymbol{\Sigma}_\perp) \\ & \geq \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) \lambda_{k-i+1}(\boldsymbol{\Sigma}) + \sum_{i=1}^{n-k} \lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) \lambda_{n-k-i+1}(\boldsymbol{\Sigma}_\perp) \\ & \geq \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) \lambda_k(\mathbf{A}) + \sum_{i=1}^{n-k} \lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) \lambda_{k+1}(\mathbf{A}), \end{aligned}$$

where we find that

$$\begin{aligned} \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) &= \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}) - \sum_{i=1}^k \lambda_i^2(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}) \\ &= \text{tr}(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}) - \text{tr}((\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U})^2) \\ &= \|\mathbf{X}^\top \mathbf{U}\|_F^2 - \|\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}\|_F^2 \end{aligned}$$

and similarly

$$\sum_{i=1}^k \lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) = -(\|\mathbf{X}_\perp^\top \mathbf{U}_\perp\|_F^2 - \|\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp\|_F^2).$$

Note that

$$\begin{aligned} \|\mathbf{X}_\perp^\top \mathbf{U}_\perp\|_F^2 &= \text{tr}(\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) = \text{tr}(\mathbf{U}_\perp \mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top) = \text{tr}((\mathbf{I} - \mathbf{U} \mathbf{U}^\top)(\mathbf{I} - \mathbf{X} \mathbf{X}^\top)) \\ &= \text{tr}(\mathbf{I} - \mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top) \\ &= n - 2k + \|\mathbf{X}^\top \mathbf{U}\|_F^2 \end{aligned}$$

and

$$\begin{aligned} & \|\mathbf{U}_\perp^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp\|_F^2 \\ &= \text{tr}((\mathbf{I} - \mathbf{U} \mathbf{U}^\top)(\mathbf{I} - \mathbf{X} \mathbf{X}^\top)(\mathbf{I} - \mathbf{U} \mathbf{U}^\top)(\mathbf{I} - \mathbf{X} \mathbf{X}^\top)) \\ &= \text{tr}((\mathbf{I} - \mathbf{U} \mathbf{U}^\top)(\mathbf{I} - \mathbf{X} \mathbf{X}^\top)(\mathbf{I} - \mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top)) \\ &= \text{tr}((\mathbf{I} - \mathbf{U} \mathbf{U}^\top)(\mathbf{I} - \mathbf{X} \mathbf{X}^\top)(\mathbf{I} + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top)) \\ &= \text{tr}(\mathbf{I} - \mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top + (\mathbf{I} - \mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top) \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top) \\ &= \text{tr}(\mathbf{I} - \mathbf{U} \mathbf{U}^\top - \mathbf{X} \mathbf{X}^\top + \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} \mathbf{X}^\top) \\ &= n - 2k + \|\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}\|_F^2. \end{aligned}$$

Therefore, we arrive at

$$\begin{aligned} & \text{tr}(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \boldsymbol{\Sigma}) + \text{tr}(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp \boldsymbol{\Sigma}_\perp) \\ & \geq \lambda_k \sum_{i=1}^k \lambda_i(\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}) + \lambda_{k+1} \sum_{i=1}^{n-k} \lambda_i(\mathbf{U}_\perp^\top \mathbf{X} \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U}_\perp) \\ & = (\lambda_k - \lambda_{k+1}) (\|\mathbf{X}^\top \mathbf{U}\|_F^2 - \|\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}\|_F^2) \\ & = \tau (\|\mathbf{X}^\top \mathbf{U}\|_F^2 - \|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}\|_F^2). \end{aligned}$$

□

**Lemma 5.5.** If  $\alpha^2 (\beta^2 + 48k\tilde{\beta}^2) < 1$ , then

$$\begin{aligned} \text{tr} (a_1(\mathbf{X}_t) a_2^{-1}(\mathbf{X}_t)) &\geq \|\mathbf{X}_t^\top \mathbf{U}\|_F^2 + 2\alpha\tau \left( \|\mathbf{X}_t^\top \mathbf{U}\|_F^2 - \|\mathbf{U}^\top \mathbf{X}_t \mathbf{X}_t^\top \mathbf{U}\|_F^2 \right) \\ &\quad - 2(1 + 2\alpha\beta) \alpha^2 \beta^2 \Theta(\mathbf{X}_t, \mathbf{U}) - k(1 + 2\alpha\beta) \alpha^2 \nu_t^2. \end{aligned}$$

*Proof.* First, we have  $a_1(\mathbf{X}) \succcurlyeq \mathbf{0}$  and  $a_2(\mathbf{X}) \succ \mathbf{0}$ . By Lemma 3 in (Shamir, 2016), we then get

$$\text{tr} (a_1(\mathbf{X}) a_2^{-1}(\mathbf{X})) \geq \text{tr} (a_1(\mathbf{X}) (2\mathbf{I} - a_2(\mathbf{X}))),$$

where if  $\alpha^2 (\beta^2 + 48k\tilde{\beta}^2) < 1$  then

$$\begin{aligned} 2\mathbf{I} - a_2(\mathbf{X}) &= \mathbf{I} - \alpha^2 \mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X} - \alpha^2 \nu_t^2 \mathbf{I} \\ &\succcurlyeq (1 - \alpha^2 \|\mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}\|_2 - \alpha^2 \nu_t^2) \mathbf{I} \succcurlyeq (1 - \alpha^2 (\beta^2 + 48k\tilde{\beta}^2)) \mathbf{I} \succ \mathbf{0}. \end{aligned}$$

Then by Lemma 2 in (Shamir, 2016) we have

$$\begin{aligned} \text{tr} (a_1(\mathbf{X}) a_2^{-1}(\mathbf{X})) &\geq \text{tr} ((a_1(\mathbf{X}) - \alpha^2 \mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) (2\mathbf{I} - a_2(\mathbf{X}))) \\ &= \text{tr} (\tilde{a}_1(\mathbf{X})) - \alpha^2 \text{tr} (\tilde{a}_1(\mathbf{X}) \mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) - \alpha^2 \nu_t^2 \text{tr} (\tilde{a}_1(\mathbf{X})), \end{aligned}$$

where  $\tilde{a}_1(\mathbf{X}) = \mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X} + 2\alpha \text{sym}(\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X})$ . We now further lower bound the right hand side of the above inequality. By Lemma 5.4., we get

$$\begin{aligned} \text{tr} (\tilde{a}_1(\mathbf{X})) &= \text{tr} (\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}) + 2\alpha \text{tr} (\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) \\ &\geq \|\mathbf{X}^\top \mathbf{U}\|_F^2 + 2\alpha\tau \left( \|\mathbf{X}^\top \mathbf{U}\|_F^2 - \|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}\|_F^2 \right). \end{aligned}$$

On the other hand, by the Cauchy-Schwarz inequality, we can obtain

$$\begin{aligned} \text{tr} (\tilde{a}_1(\mathbf{X})) &= \text{tr} (\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}) + 2\alpha \text{tr} (\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) \\ &\leq \|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top\|_F \|\mathbf{X}\|_F + 2\alpha \|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A}\|_F \|\mathbf{X}\|_F \\ &\leq \|\mathbf{X}^\top\|_F \|\mathbf{U} \mathbf{U}^\top\|_2 \|\mathbf{X}\|_F + 2\alpha \|\mathbf{X}^\top\|_F \|\mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A}\|_2 \|\mathbf{X}\|_F \\ &\leq (1 + 2\alpha\beta) \|\mathbf{X}\|_F^2 = (1 + 2\alpha\beta) k. \end{aligned}$$

In addition,

$$\begin{aligned} &\text{tr} (\tilde{a}_1(\mathbf{X}) \mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}) \\ &\leq \|\tilde{a}_1(\mathbf{X}) \mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp\|_F \|\mathbf{X}_\perp^\top \mathbf{A} \mathbf{X}\|_F \leq \|\tilde{a}_1(\mathbf{X})\|_2 \|\mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp\|_F^2 \\ &\leq (\|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}\|_2 + 2\alpha \|\mathbf{X}^\top \mathbf{U} \mathbf{U}^\top \mathbf{X}_\perp \mathbf{X}_\perp^\top \mathbf{A}\|_2) \|\mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp\|_F^2 \\ &\leq (1 + 2\alpha\beta) \|\mathbf{X}^\top \mathbf{A} \mathbf{X}_\perp\|_F^2 = (1 + 2\alpha) \|\mathbf{X}^\top (\mathbf{U} \Sigma \mathbf{U}^\top + \mathbf{U}_\perp \Sigma_\perp \mathbf{U}_\perp^\top) \mathbf{X}_\perp\|_F^2 \\ &\leq (1 + 2\alpha\beta) (\|\mathbf{X}^\top \mathbf{U} \Sigma \mathbf{U}^\top \mathbf{X}_\perp\|_F + \|\mathbf{X}^\top \mathbf{U}_\perp \Sigma_\perp \mathbf{U}_\perp^\top \mathbf{X}_\perp\|_F)^2 \\ &\leq (1 + 2\alpha\beta) (\|\mathbf{X}^\top \mathbf{U} \Sigma\|_2 \|\mathbf{U}^\top \mathbf{X}_\perp\|_F + \|\mathbf{X}^\top \mathbf{U}_\perp\|_F \|\Sigma_\perp \mathbf{U}_\perp^\top \mathbf{X}_\perp\|_2)^2 \\ &\leq (1 + 2\alpha\beta) (\|\mathbf{A}\|_2 \|\mathbf{U}^\top \mathbf{X}_\perp\|_F + \|\mathbf{X}^\top \mathbf{U}_\perp\|_F \|\mathbf{A}\|_2)^2 \\ &\leq \beta^2 (1 + 2\alpha\beta) (\|\mathbf{U}^\top \mathbf{X}_\perp\|_F + \|\mathbf{X}^\top \mathbf{U}_\perp\|_F)^2 = 2\beta^2 (1 + 2\alpha\beta) \Theta(\mathbf{X}, \mathbf{U}). \end{aligned}$$

Therefore, we arrive at

$$\begin{aligned} \text{tr} (a_1(\mathbf{X}) a_2^{-1}(\mathbf{X})) &\geq \|\mathbf{X}^\top \mathbf{U}\|_F^2 + 2\alpha\tau \left( \|\mathbf{X}^\top \mathbf{U}\|_F^2 - \|\mathbf{U}^\top \mathbf{X} \mathbf{X}^\top \mathbf{U}\|_F^2 \right) \\ &\quad - 2(1 + 2\alpha\beta) \alpha^2 \beta^2 \Theta(\mathbf{X}, \mathbf{U}) - k(1 + 2\alpha\beta) \alpha^2 \nu_t^2. \end{aligned}$$

□

**Lemma 5.6.** For any  $\iota \in (0, 1)$ , if  $\alpha$  and  $m$  satisfy that  $\alpha < 2\xi\tau/\gamma$ ,  $\min\{2(1 + \alpha\beta)\alpha\mu, \alpha^2(\beta^2 + 48k\tilde{\beta}^2)\} < 1$  and  $\Theta_0 + km\rho + \theta\sqrt{2m\log(1/\iota)} < \delta$ , then  $\Theta_t < \delta$  holds for all  $t = 1, 2, \dots, m$  with probability at least  $1 - \iota$ , where  $\theta = \frac{4k(\beta+4\tilde{\beta})\alpha}{1-(\beta+4\tilde{\beta})\alpha} + k\rho$ .

*Proof.* Consider the stochastic process  $\{\Theta_0, \Theta_1, \dots, \Theta_m\}$  and the filtration  $\mathcal{F} = \{\mathcal{F}_t\}$  about random draws  $y_t$ . Since  $\mathbb{E}[\Theta_{t+1}|\mathbf{X}_t] \leq (1 - \alpha(2\xi\tau - \gamma\alpha))\Theta_t + \rho\Theta_{s-1}$ , we have  $\mathbb{E}[\Theta_{t+1}|X^{(t)}] \leq \Theta_t + k\rho$  if  $\alpha < 2\xi\tau/\gamma$ . Define  $\Psi_t = \Theta_t - k\rho t$  for  $t = 0, 1, 2, \dots, m$ . Note that  $\{\Psi_t : t = 0, 1, 2, \dots, m\}$  is a finite sequence of random variables and thus the natural continuation can be applied to get an infinite sequence such that  $|\Psi_t| \leq \Theta_t + k\rho m \leq k(1 + \rho m)$  for all  $t$  including  $t > m$ . In addition,

$$\mathbb{E}[\Psi_t|\mathbf{X}_{t-1}] = \mathbb{E}[\Theta_t|\mathbf{X}_{t-1}] - k\rho t \leq \Theta_{t-1} + k\rho - k\rho t = \Theta_{t-1} - k\rho(t-1) = \Psi_{t-1}.$$

$\{\Psi_t\}$  thus is a super-martingale. Next, we prove that it has bounded differences. Note that the intermediate update in the tangent space is  $\mathbf{X}_t + \alpha_{t+1}\text{Grad}f(\mathbf{X}_t) + \alpha_{t+1}\mathbf{W}_t$ , for which we have

$$\|\text{Grad}f(\mathbf{X}_t) + \mathbf{W}_t\|_2 \leq \|\text{Grad}f(\mathbf{X}_t)\|_2 + \|\mathbf{W}_t\|_2 \leq \beta + 4\tilde{\beta}$$

using Lemma 5.1. Based on the proof of Lemma 9 in (Shamir, 2016), we have that

$$\left| \|\mathbf{X}_{t+1}^\top \mathbf{U}\|_F^2 - \|\mathbf{X}_t^\top \mathbf{U}\|_F^2 \right| \leq \frac{4k(\beta + 4\tilde{\beta})\alpha}{1 - (\beta + 4\tilde{\beta})\alpha}.$$

Thus, we get the bounded difference

$$|\Psi_{t+1} - \Psi_t| \leq \left| \|\mathbf{X}_{t+1}^\top \mathbf{U}\|_F^2 - \|\mathbf{X}_t^\top \mathbf{U}\|_F^2 \right| + k\rho \leq \frac{4k(\beta + 4\tilde{\beta})\alpha}{1 - (\beta + 4\tilde{\beta})\alpha} + k\rho \triangleq \theta.$$

We now are able to apply the Azuma-Hoeffding inequality, i.e., for any  $0 \leq t \leq m$  and  $a > 0$

$$P(\Psi_t - \Psi_0 \geq a) \leq \exp\left\{-\frac{a^2}{2\sum_{s=1}^t \theta^2}\right\} = \exp\left\{-\frac{a^2}{2t\theta^2}\right\} \leq \exp\left\{-\frac{a^2}{2m\theta^2}\right\} \triangleq \iota \in (0, 1).$$

Solving  $\iota = \exp\left\{-\frac{a^2}{2m\theta^2}\right\}$  with respect to  $a$  yields  $a = \theta\sqrt{2m\log(1/\iota)}$ . Therefore, we get that  $\Psi_t - \Psi_0 < a$ , i.e.,

$$\Theta_t \leq \Theta_0 + k\rho t + a \leq \Theta_0 + km\rho + \theta\sqrt{2m\log(1/\iota)} < \delta$$

for all  $t = 1, 2, \dots, m$ , with probability at least  $1 - \iota$ . □