# A  Proof of the Main Theory

In this section, we present the proofs of our main theorems. Let us first recall the notations used in Algorithm 1. $\mathbf{v}_t^s$ and $\mathbf{U}_t^s$ are the semi-stochastic gradient and Hessian defined in (3.1) and (3.2) respectively. $\mathbf{x}_t^s$'s are the iterates and $\widehat{\mathbf{x}}^s$'s are the reference points used in Algorithm 1. $b_g$ and $b_h$ are the batch sizes of semi-stochastic gradient and Hessian. $S$ and $T$ are the number of epochs and epoch length of Algorithm 1. We set $M_{s,t} := M = C_M\rho$ as suggested by Theorems 4.7 and 5.3, where $C_M > 0$ is a constant. $\mathbf{h}_t^s$ is the exact minimizer of $m_t^s(\mathbf{h})$, where $m_t^s(\mathbf{h})$ is defined in (3.3). $\widetilde{\mathbf{h}}_t^s$ is the inexact minimizer defined in (5.1).

In order to prove Theorems 4.7 and 5.3, we first lay down the following useful technical lemmas.

**Lemma A.1.** (Nesterov & Polyak, 2006) We have the following basic results:

$$\mathbf{v}_t^s + \mathbf{U}_t^s \mathbf{h}_t^s + \frac{M}{2}\|\mathbf{h}_t^s\|_2 \mathbf{h}_t^s = 0, \tag{A.1}$$

$$\mathbf{U}_t^s + \frac{M}{2}\|\mathbf{h}_t^s\|_2 \mathbf{I} \succeq 0, \tag{A.2}$$

$$\langle \mathbf{v}_t^s, \mathbf{h}_t^s \rangle + \frac{1}{2}\langle \mathbf{U}_t^s \mathbf{h}_t^s, \mathbf{h}_t^s \rangle + \frac{M}{6}\|\mathbf{h}_t^s\|_2^3 \leq -\frac{M}{12}\|\mathbf{h}_t^s\|_2^3. \tag{A.3}$$

Next we have two lemmas which we use to control the variance of $\mathbf{v}_t^s$ and $\mathbf{U}_t^s$. They play important roles in our proof:

**Lemma A.2.** For the semi-stochastic gradient $\mathbf{v}_t^s$ defined in (3.1), we have

$$\mathbb{E}_{i_t}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} \leq \frac{\rho^{3/2}}{b_g^{3/4}}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3,$$

where $\mathbb{E}_{i_t}$ is the expectation over all $i_t \in I_g$.

**Lemma A.3.** Let $\mathbf{U}_t^s$ be the semi-stochastic Hessian defined in (3.2). If the batch size satisfy $b_h \geq 400 \log d$, then we have

$$\mathbb{E}_{j_t}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3 \leq 1200\rho^3\left(\frac{\log d}{b_h}\right)^{3/2}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3,$$

where $\mathbb{E}_{j_t}$ is the expectation over all $j_t \in I_h$.

**Lemma A.4.** For any $\mathbf{h} \in \mathbb{R}^d$, we have

$$\langle \nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s, \mathbf{h} \rangle \leq \frac{M}{27}\|\mathbf{h}\|_2^3 + \frac{2\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{M^{1/2}},$$

$$\langle (\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s)\mathbf{h}, \mathbf{h} \rangle \leq \frac{2M}{27}\|\mathbf{h}\|_2^3 + \frac{27}{M^2}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3.$$

**Lemma A.5.** For any $\mathbf{h} \in \mathbb{R}^d$, if $C_M \geq 100$, then we have

$$\mu(\mathbf{x}_t^s + \mathbf{h}) \leq 9C_M^{3/2}\Big[M^{3/2}\|\mathbf{h}\|_2^3 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + M^{-3/2}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3$$
$$+ \|\nabla m_t^s(\mathbf{h})\|_2^{3/2} + M^{3/2}\big|\|\mathbf{h}\|_2 - \|\mathbf{h}_t^s\|_2\big|^3\Big].$$

**Lemma A.6.** If $T \geq 1$, then for any $\mathbf{h} \in \mathbb{R}^d$, we have

$$\|\mathbf{x}_t^s + \mathbf{h} - \widehat{\mathbf{x}}^s\|_2^3 \leq 2T^2\|\mathbf{h}\|_2^3 + (1 + 3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3. \tag{A.4}$$

**Lemma A.7.** We define constant series $c_t, 0 \leq t \leq T$ as the following: $c_T = 0$, $c_t = c_{t+1}(1 + 3/T) + M(500T^3)^{-1}$, $0 \leq t \leq T - 1$. Then we have for any $1 \leq t \leq T$,

$$M/24 - 2c_t T^2 \geq 0. \tag{A.5}$$

## A.1 Proof of Theorem 4.7

*Proof of Theorem 4.7.* We first upper bound $F(\mathbf{x}_{t+1}^s)$ as follows

$$
\begin{aligned}
F(\mathbf{x}_{t+1}^s) &\leq F(\mathbf{x}_t^s) + \langle \nabla F(\mathbf{x}_t^s), \mathbf{h}_t^s \rangle + \frac{1}{2}\langle \nabla^2 F(\mathbf{x}_t^s)\mathbf{h}_t^s, \mathbf{h}_t^s \rangle + \frac{\rho}{6}\|\mathbf{h}_t^s\|_2^3 \\
&= F(\mathbf{x}_t^s) + \langle \mathbf{v}_t^s, \mathbf{h}_t^s \rangle + \frac{1}{2}\langle \mathbf{U}_t^s\mathbf{h}_t^s, \mathbf{h}_t^s \rangle + \frac{M}{6}\|\mathbf{h}_t^s\|_2^3 + \langle \nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s, \mathbf{h}_t^s \rangle \\
&\quad + \frac{1}{2}\langle (\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s)\mathbf{h}_t^s, \mathbf{h}_t^s \rangle - \frac{M-\rho}{6}\|\mathbf{h}_t^s\|_2^3 \\
&\leq F(\mathbf{x}_t^s) - \frac{M}{12}\|\mathbf{h}_t^s\|_2^3 + \left( \frac{M}{27}\|\mathbf{h}_t^s\|_2^3 + \frac{2\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{M^{1/2}} \right) \\
&\quad + \frac{1}{2}\left( \frac{2M}{27}\|\mathbf{h}_t^s\|_2^3 + \frac{27}{M^2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 \right) - \frac{M-\rho}{6}\|\mathbf{h}_t^s\|_2^3 \\
&\leq F(\mathbf{x}_t^s) - \frac{M}{12}\|\mathbf{h}_t^s\|_2^3 + \frac{2}{M^{1/2}}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + \frac{27}{M^2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3, \quad (A.6)
\end{aligned}
$$

where the first inequality follows from Lemma 4.2 and the second inequality holds due to Lemmas A.1 and A.4. We define

$$
R_t^s = \mathbb{E}\big[ F(\mathbf{x}_t^s) + c_t\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \big], \quad (A.7)
$$

where $c_t$ is defined in the Lemma A.7. Then by Lemma A.6, we have

$$
c_{t+1}\|\mathbf{x}_{t+1}^s - \widehat{\mathbf{x}}^s\|_2^3 \leq 2c_{t+1}T^2\|\mathbf{h}_t^s\|_2^3 + c_{t+1}(1 + 3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3. \quad (A.8)
$$

Applying Lemma A.5 with $\mathbf{h} = \mathbf{h}_t^s$, we have

$$
\begin{aligned}
\big(240C_M^2\rho^{1/2}\big)^{-1}\mu(\mathbf{x}_{t+1}^s) &\leq \frac{M}{24}\|\mathbf{h}_t^s\|_2^3 + \frac{\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{24M^{1/2}} + \frac{\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3}{24M^2} \\
&\quad + \frac{\|\nabla m_t^s(\mathbf{h}_t^s)\|_2^{3/2}}{24M^{1/2}} + \frac{M}{24}\big| \|\mathbf{h}_t^s\|_2 - \|\mathbf{h}_t^s\|_2 \big|^3 \\
&= \frac{M}{24}\|\mathbf{h}_t^s\|_2^3 + \frac{\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{24M^{1/2}} + \frac{\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3}{24M^2}, \quad (A.9)
\end{aligned}
$$

where the equality is due to Lemma A.1.

Adding (A.6) with (A.8) and (A.9) and taking total expectation, we have

$$
\begin{aligned}
&R_{t+1}^s + \big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_{t+1}^s) \\
&= \mathbb{E}\Big[ F(\mathbf{x}_{t+1}^s) + c_{t+1}\|\mathbf{x}_{t+1}^s - \widehat{\mathbf{x}}^s\|_2^3 + \big(240C_M^2\rho^{1/2}\big)^{-1}\mu(\mathbf{x}_{t+1}^s) \Big] \\
&\leq \mathbb{E}\Big[ F(\mathbf{x}_t^s) + c_{t+1}(1+3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 - \|\mathbf{h}_t^s\|_2^3\big(M/24 - 2c_{t+1}T^2\big) \Big] \\
&\quad + \mathbb{E}\Big[ 3M^{-1/2}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + 28M^{-2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 \Big] \\
&\leq \mathbb{E}\Big[ F(\mathbf{x}_t^s) + c_{t+1}(1+3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \Big] + \mathbb{E}\Big[ 3M^{-1/2}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + 28M^{-2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 \Big], \quad (A.10)
\end{aligned}
$$

where the third inequality holds due to Lemma A.7. To further bound (A.10), we have

$$
\frac{3}{M^{1/2}}\mathbb{E}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} \leq \frac{3\rho^{3/2}}{M^{1/2}b_g^{3/4}}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \leq \frac{M}{1000T^3}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3, \quad (A.11)
$$

where the first inequality holds due to Lemma A.2, the second inequality holds due to $M \geq 100\rho$ and $b_g \geq 5T^4$ from the condition of Theorem 4.7. We also have

$$
\frac{28}{M^2}\mathbb{E}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 \leq \frac{28 \times 15000\rho^3}{M^2(b_h/\log d)^{3/2}}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \leq \frac{M}{1000T^3}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3, \quad (A.12)
$$

where the first inequality holds due to Lemma A.3, where we have $b_h \geq 100T^2 \log d \geq 400 \log d$, and the second inequality holds due to $M \geq 100\rho$ and $b_h \geq 100T^2 \log d$ from the assumption of Theorem 4.7. Thus, submitting (A.11) and (A.12) into (A.10), we have

$$
\begin{aligned}
R_{t+1}^s + \left(240C_M^2\rho^{1/2}\right)^{-1}\mathbb{E}\mu(\mathbf{x}_{t+1}^s) &\leq \mathbb{E}\left[F(\mathbf{x}_t^s) + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\left(c_{t+1}(1+3/T) + \frac{M}{500T^3}\right)\right] \\
&= \mathbb{E}\left[F(\mathbf{x}_t^s) + c_t\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\right] \\
&= R_t^s,
\end{aligned}
\tag{A.13}
$$

where the first equality holds due to the definition of $c_t$ in Lemma A.7. Telescoping (A.13) from $t=0$ to $T-1$, we have

$$
R_0^s - R_T^s \geq \sum_{t=1}^{T}\left(240C_M^2\rho^{1/2}\right)^{-1}\mathbb{E}\mu(\mathbf{x}_t^s).
$$

By the definition of $c_T$ in Lemma A.7, we have $c_T = 0$, then $R_T^s = \mathbb{E}\left[F(\mathbf{x}_T^s) + c_T\|\mathbf{x}_T^s - \widehat{\mathbf{x}}^s\|_2^3\right] = \mathbb{E}F(\widehat{\mathbf{x}}^{s+1})$; meanwhile by the definition of $\mathbf{x}_0^s$, we have $\mathbf{x}_0^s = \widehat{\mathbf{x}}^s$. Thus we have $R_0^s = \mathbb{E}\left[F(\mathbf{x}_0^s) + c_0\|\mathbf{x}_0^s - \widehat{\mathbf{x}}^s\|_2^3\right] = \mathbb{E}F(\widehat{\mathbf{x}}^s)$, which implies

$$
\mathbb{E}F(\widehat{\mathbf{x}}^s) - \mathbb{E}F(\widehat{\mathbf{x}}^{s+1}) = R_0^s - R_T^s \geq \left(240C_M^2\rho^{1/2}\right)^{-1}\sum_{t=1}^{T}\mathbb{E}\mu(\mathbf{x}_t^s).
\tag{A.14}
$$

Finally, telescoping (A.14) from $s=1$ to $S$ yields

$$
\Delta_F \geq \sum_{s=1}^{S}\mathbb{E}F(\widehat{\mathbf{x}}^s) - \mathbb{E}F(\widehat{\mathbf{x}}^{s+1}) \geq \left(240C_M^2\rho^{1/2}\right)^{-1}\sum_{s=1}^{S}\sum_{t=1}^{T}\mathbb{E}\mu(\mathbf{x}_t^s).
$$

By the definition about choice of $\mathbf{x}_{\text{out}}$, we complete the proof. $\square$

## A.2  Proof of Corollary 4.10

*Proof.* We can verify that the parameter setting in Corollary 4.10 satisfies the requirement of Theorem 4.7. Thus, submitting the choice of parameters into Theorem 4.7, the output of Algorithm 1 $\mathbf{x}_{\text{out}}$ satisfies that

$$
\mathbb{E}[\mu(\mathbf{x}_{\text{out}})] \leq \frac{240C_M^2\rho^{1/2}\Delta_F}{ST} \leq \epsilon^{3/2},
\tag{A.15}
$$

which indeed implies that $\mathbf{x}_{\text{out}}$ is an $(\epsilon, \sqrt{\rho\epsilon})$-approximate local minimum. Next we calculate how many SO calls and CSO calls are needed. Algorithm 1 needs to calculate full gradient $\mathbf{g}_s$ and full Hessian $\mathbf{H}_s$ at the beginning of each epoch, with $n$ SO calls. In each epoch, Algorithm 1 needs to calculate $\mathbf{v}_t^s$ and $\mathbf{U}_t^s$ with $b_g + b_h$ SO calls at each iteration. Thus, the total amount of SO calls is

$$
\begin{aligned}
Sn + (ST)(b_g + b_h) &\leq n + C_1\Delta_F\rho^{1/2}n^{4/5}\epsilon^{-3/2} + C_1\Delta_F\rho^{1/2}\epsilon^{-3/2}(5n^{4/5} + 1000n^{2/5}\log d) \\
&= \widetilde{O}\left(n + \frac{\Delta_F\sqrt{\rho}n^{4/5}}{\epsilon^{3/2}}\right),
\end{aligned}
$$

where $C_1 = 240C_M^2$. For the CSO calls, Algorithm 1 needs to solve cubic subproblem at each single iteration. Thus, the total amount of CSO calls is

$$
ST \leq C_1\Delta_F\rho^{1/2}\epsilon^{-3/2} = O\left(\frac{\Delta_F\sqrt{\rho}}{\epsilon^{3/2}}\right).
$$

$\square$

## A.3  Proof of Theorem 5.3

*Proof of Theorem 5.3.* Similar to (A.6) in the proof of Theorem 4.7, we have

$$
\begin{aligned}
F(\mathbf{x}_{t+1}^s) &\leq F(\mathbf{x}_t^s) + \langle \nabla F(\mathbf{x}_t^s), \widetilde{\mathbf{h}}_t^s \rangle + \frac{1}{2}\langle \nabla^2 F(\mathbf{x}_t^s)\widetilde{\mathbf{h}}_t^s, \widetilde{\mathbf{h}}_t^s \rangle + \frac{\rho}{6}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 \\
&= F(\mathbf{x}_t^s) + \langle \mathbf{v}_t^s, \widetilde{\mathbf{h}}_t^s \rangle + \frac{1}{2}\langle \mathbf{U}_t^s\widetilde{\mathbf{h}}_t^s, \widetilde{\mathbf{h}}_t^s \rangle + \frac{M}{6}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \langle \nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s, \widetilde{\mathbf{h}}_t^s \rangle \\
&\quad + \frac{1}{2}\langle (\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s)\widetilde{\mathbf{h}}_t^s, \widetilde{\mathbf{h}}_t^s \rangle - \frac{M-\rho}{6}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 \\
&\leq F(\mathbf{x}_t^s) - \frac{M}{12}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \delta + \bigg(\frac{M}{27}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \frac{2\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{M^{1/2}}\bigg) \\
&\quad + \frac{1}{2}\bigg(\frac{2M}{27}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \frac{27}{M^2}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3\bigg) - \frac{M-\rho}{6}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 \\
&\leq F(\mathbf{x}_t^s) - \frac{M}{12}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \frac{2}{M^{1/2}}\big\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\big\|_2^{3/2} + \frac{27}{M^2}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3 + \delta, \quad \text{(A.16)}
\end{aligned}
$$

where the second inequality holds because $\widetilde{\mathbf{h}}_t^s$ is an inexact solver satisfying Condition 5.1. By Lemma A.6 with $\mathbf{h} = \widetilde{\mathbf{h}}_t^s$, we have

$$
c_{t+1}\|\mathbf{x}_{t+1}^s - \widehat{\mathbf{x}}^s\|_2^3 = c_{t+1}\big\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s + \widetilde{\mathbf{h}}_t^s\big\|_2^3 \leq 2c_{t+1}T^2\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + c_{t+1}(1 + 3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3. \quad \text{(A.17)}
$$

By Lemma A.5, we also have

$$
\begin{aligned}
&\big(240C_M^2\rho^{1/2}\big)^{-1}\mu(\mathbf{x}_{t+1}^s) \\
&= \big(240C_M^2\rho^{1/2}\big)^{-1}\mu\big(\mathbf{x}_t^s + \widetilde{\mathbf{h}}_t^s\big) \\
&\leq \frac{M}{24}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \frac{\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{24M^{1/2}} + \frac{\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3}{24M^2} + \frac{\big\|\nabla m_t^s\big(\widetilde{\mathbf{h}}_t^s\big)\big\|_2^{3/2}}{24M^{1/2}} + \frac{M\big|\big\|\widetilde{\mathbf{h}}_t^s\big\|_2 - \|\mathbf{h}_t^s\|_2\big|^3}{24}, \quad \text{(A.18)}
\end{aligned}
$$

Since $\widetilde{\mathbf{h}}_t^s$ is an inexact solver satisfying Condition 5.1, we have

$$
\frac{\big\|\nabla m_t^s\big(\widetilde{\mathbf{h}}_t^s\big)\big\|_2^{3/2}}{24M^{1/2}} + \frac{M\big|\big\|\widetilde{\mathbf{h}}_t^s\big\|_2 - \|\mathbf{h}_t^s\|_2\big|^3}{24} \leq \frac{\delta}{24} + \frac{\delta}{24} < \delta. \quad \text{(A.19)}
$$

Submitting (A.19) into (A.18), we have

$$
\big(240C_M^2\rho^{1/2}\big)^{-1}\mu(\mathbf{x}_{t+1}^s) \leq \frac{M}{24}\big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3 + \frac{\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{24M^{1/2}} + \frac{\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3}{24M^2} + \delta. \quad \text{(A.20)}
$$

Then adding (A.16), (A.17) and (A.20) up, we have

$$
\begin{aligned}
&R_{t+1}^s + \big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_{t+1}^s) \\
&= \mathbb{E}\Big[F(\mathbf{x}_{t+1}^s) + c_{t+1}\|\mathbf{x}_{t+1}^s - \widehat{\mathbf{x}}^s\|_2^3 + \big(240C_M^2\rho^{1/2}\big)^{-1}\mu(\mathbf{x}_{t+1}^s)\Big] \\
&\leq \mathbb{E}\Big[F(\mathbf{x}_t^s) + c_{t+1}(1 + 3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 - \big\|\widetilde{\mathbf{h}}_t^s\big\|_2^3\big(M/24 - 2c_{t+1}T^2\big)\Big] \\
&\quad + \mathbb{E}\Big[\frac{3}{M^{1/2}}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + \frac{28}{M^2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3\Big] + 2\delta \\
&\leq \mathbb{E}\Big[F(\mathbf{x}_t^s) + c_{t+1}(1 + 3/T)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\Big] + \mathbb{E}\Big[\frac{3}{M^{1/2}}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + \frac{28}{M^2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3\Big] + 2\delta. \quad \text{(A.21)}
\end{aligned}
$$

Since the parameter setting is the same as Theorem 4.7, by (A.11) and (A.12), we have

$$
\frac{3}{M^{1/2}}\mathbb{E}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} \leq \frac{M}{1000T^3}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3, \quad \text{(A.22)}
$$

and

$$\frac{28}{M^2}\mathbb{E}\big\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\big\|_2^3 \leq \frac{M}{1000T^3}\mathbb{E}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3. \tag{A.23}$$

Submitting (A.22) and (A.23) into (A.21) yields

$$
\begin{aligned}
R_{t+1}^s + \big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_{t+1}^s) &\leq \mathbb{E}\bigg[F(\mathbf{x}_t^s) + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\Big(c_{t+1}(1+3/T) + \frac{M}{500T^3}\Big)\bigg] + 2\delta \\
&= \mathbb{E}\big[F(\mathbf{x}_t^s) + c_t\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\big] + 2\delta \\
&= R_t^s + 2\delta,
\end{aligned}
\tag{A.24}
$$

where the first equality holds due to the definition of $c_t$ in Lemma A.7. Telescoping (A.13) from $t = 0$ to $T - 1$, we have

$$R_0^s - R_T^s \geq \sum_{t=1}^{T}\Big[\big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_t^s) - 2\delta\Big].$$

By the definition of $c_T$ in Lemma A.7, we have $c_T = 0$, then $R_T^s = \mathbb{E}\big[F(\mathbf{x}_T^s) + c_T\|\mathbf{x}_T^s - \widehat{\mathbf{x}}^s\|_2^3\big] = \mathbb{E}F(\widehat{\mathbf{x}}^{s+1})$; meanwhile by the definition of $\mathbf{x}_0^s$, we have $\mathbf{x}_0^s = \widehat{\mathbf{x}}^s$. Thus we have $R_0^s = \mathbb{E}\big[F(\mathbf{x}_0^s) + c_0\|\mathbf{x}_0^s - \widehat{\mathbf{x}}^s\|_2^3\big] = \mathbb{E}F(\widehat{\mathbf{x}}^s)$, which further implies

$$\mathbb{E}F(\widehat{\mathbf{x}}^s) - \mathbb{E}F(\widehat{\mathbf{x}}^{s+1}) = R_0^s - R_T^s \geq \sum_{t=1}^{T}\Big[\big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_t^s) - 2\delta\Big]. \tag{A.25}$$

Finally, telescoping (A.25) from $s = 1$ to $S$, we obtain

$$\Delta_F \geq \sum_{s=1}^{S}\mathbb{E}F(\widehat{\mathbf{x}}^s) - \mathbb{E}F(\widehat{\mathbf{x}}^{s+1})\sum_{s=1}^{S}\sum_{t=1}^{T}\Big[\big(240C_M^2\rho^{1/2}\big)^{-1}\mathbb{E}\mu(\mathbf{x}_t^s) - 2\delta\Big].$$

By the definition about choice of $\mathbf{x}_{\text{out}}$, we finish the proof.

$\square$

## A.4 Proof of Corollary 5.5

*Proof of Corollary 5.5.* We can verify that under the parameter choice in Corollary 5.5,

$$\mathbb{E}[\mu(\mathbf{x}_{\text{out}})] \leq \frac{240C_M^2\rho^{1/2}\Delta_F}{ST} + 480C_M^2\rho^{1/2}\delta \leq \epsilon^{3/2}/2 + \epsilon^{3/2}/2 = \epsilon^{3/2}. \tag{A.26}$$

Thus, $\mathbf{x}_{\text{out}}$ is an $(\epsilon, \sqrt{\rho\epsilon})$-approximate local minimum. By the proof of Corollary 4.10, the total amount of SO calls is

$$
\begin{aligned}
Sn + (ST)(b_g + b_h) &\leq n + C_1\Delta_F\rho^{1/2}n^{4/5}\epsilon^{-3/2} + C_1\Delta_F\rho^{1/2}\epsilon^{-3/2}(5n^{4/5} + 1000n^{2/5}\log d) \\
&= \widetilde{O}\Big(n + \frac{\Delta_F\sqrt{\rho}n^{4/5}}{\epsilon^{3/2}}\Big),
\end{aligned}
$$

where $C_1 = 480C_M^2$. For the CSO calls, Algorithm 1 needs to solve cubic subproblem at each single iteration. Thus, the total amount of CSO calls is

$$ST \leq C_1\Delta_F\rho^{1/2}\epsilon^{-3/2} = O\Big(\frac{\Delta_F\sqrt{\rho}}{\epsilon^{3/2}}\Big).$$

$\square$

# B Proof of Technical Lemmas

Now we prove the technical lemmas used in Section A.

## B.1  Proof of Lemma A.1

The result of Lemma A.1 is typical in the literature of cubic regularization (Nesterov & Polyak, 2006; Cartis et al., 2011a;b), but no exactly the same result has been shown in any formal way. Thus we present the proof here for self-containedness.

*Proof of Lemma A.1.* For simplicity, we let $\mathbf{g} = \mathbf{v}_t^s, \mathbf{H} = \mathbf{U}_t^s, \theta = M_t$ and $\mathbf{h}_{\text{opt}} = \mathbf{h}_t^s$. Then we need to prove

$$\mathbf{g} + \mathbf{H}\mathbf{h}_{\text{opt}} + \frac{\theta}{2}\|\mathbf{h}_{\text{opt}}\|_2\mathbf{h}_{\text{opt}} = \mathbf{0}, \tag{B.1}$$

$$\mathbf{H} + \frac{\theta}{2}\|\mathbf{h}_{\text{opt}}\|_2\mathbf{I} \succeq \mathbf{0}, \tag{B.2}$$

$$\langle \mathbf{g}, \mathbf{h}_{\text{opt}}\rangle + \frac{1}{2}\langle \mathbf{H}\mathbf{h}_{\text{opt}}, \mathbf{h}_{\text{opt}}\rangle + \frac{\theta}{6}\|\mathbf{h}_{\text{opt}}\|_2^3 \leq -\frac{\theta}{12}\|\mathbf{h}_{\text{opt}}\|_2^3. \tag{B.3}$$

Let $\lambda = \theta\|\mathbf{h}_{\text{opt}}\|_2/2$. Note that $\mathbf{h}_{\text{opt}} = \arg\min m(\mathbf{h})$, then the necessary condition $\nabla m(\mathbf{h}_{\text{opt}}) = \mathbf{0}$ and $\nabla^2 m(\mathbf{h}_{\text{opt}}) \succeq \mathbf{0}$ can be written as

$$\nabla m(\mathbf{h}_{\text{opt}}) = \mathbf{g} + \mathbf{H}\mathbf{h}_{\text{opt}} + \lambda\mathbf{h}_{\text{opt}} = \mathbf{0}, \tag{B.4}$$

$$\mathbf{w}^\top \nabla^2 m(\mathbf{h}_{\text{opt}})\mathbf{w} = \mathbf{w}^\top\left(\mathbf{H} + \lambda\mathbf{I} + \lambda\left(\frac{\mathbf{h}_{\text{opt}}}{\|\mathbf{h}_{\text{opt}}\|_2}\right)\left(\frac{\mathbf{h}_{\text{opt}}}{\|\mathbf{h}_{\text{opt}}\|_2}\right)^\top\right)\mathbf{w} \geq 0, \forall \mathbf{w} \in \mathbb{R}^d. \tag{B.5}$$

Apparently, (B.4) directly implies (B.1). To prove (B.2), we adapt the proof of Lemma 5.1 in Agarwal et al. (2017). Note that if $\langle \mathbf{w}, \mathbf{h}_{\text{opt}}\rangle = 0$, then (B.5) directly implies (B.2). So we only need to focus on the case that $\langle \mathbf{w}, \mathbf{h}_{\text{opt}}\rangle \neq 0$.

Since $\langle \mathbf{w}, \mathbf{h}_{\text{opt}}\rangle \neq 0$, there exists $\eta \neq 0$ such that $\|\mathbf{h}_{\text{opt}} + \eta\mathbf{w}\|_2 = \|\mathbf{h}_{\text{opt}}\|_2$. (In fact, we can find $\eta = -2\langle \mathbf{w}, \mathbf{h}_{\text{opt}}\rangle/\|\mathbf{w}\|_2^2$ satisfies the requirement). Next we will take a close look at the difference $m(\mathbf{h}_{\text{opt}} + \eta\mathbf{w}) - m(\mathbf{h}_{\text{opt}})$. On one hand, we have

$$\begin{aligned}
m(\mathbf{h}_{\text{opt}} + \eta\mathbf{w}) - m(\mathbf{h}_{\text{opt}}) &= \mathbf{g}^\top[(\mathbf{h}_{\text{opt}} + \eta\mathbf{w}) - \mathbf{h}_{\text{opt}}] + \frac{(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})^\top\mathbf{H}(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})}{2} - \frac{\mathbf{h}_{\text{opt}}^\top\mathbf{H}\mathbf{h}_{\text{opt}}}{2} \\
&= -[(\mathbf{h}_{\text{opt}} + \eta\mathbf{w}) - \mathbf{h}_{\text{opt}}]^\top(\mathbf{H} + \lambda\mathbf{I})\mathbf{h}_{\text{opt}} + \frac{(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})^\top\mathbf{H}(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})}{2} - \frac{\mathbf{h}_{\text{opt}}^\top\mathbf{H}\mathbf{h}_{\text{opt}}}{2} \quad\text{(B.6)} \\
&= \frac{\lambda\eta^2}{2}\|\mathbf{w}\|_2^2 + [\mathbf{h}_{\text{opt}} - (\mathbf{h}_{\text{opt}} + \eta\mathbf{w})]^\top\mathbf{H}\mathbf{h}_{\text{opt}} + \frac{(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})^\top\mathbf{H}(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})}{2} - \frac{\mathbf{h}_{\text{opt}}^\top\mathbf{H}\mathbf{h}_{\text{opt}}}{2} \\
&\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\text{(B.7)} \\
&= \frac{\lambda\eta^2}{2}\|\mathbf{w}\|_2^2 + \frac{\mathbf{h}_{\text{opt}}^\top\mathbf{H}\mathbf{h}_{\text{opt}}}{2} - (\mathbf{h}_{\text{opt}} + \eta\mathbf{w})^\top\mathbf{H}\mathbf{h}_{\text{opt}} + \frac{(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})^\top\mathbf{H}(\mathbf{h}_{\text{opt}} + \eta\mathbf{w})}{2} \\
&= \frac{\lambda\eta^2}{2}\|\mathbf{w}\|_2^2 + \frac{\eta^2}{2}\mathbf{w}^\top\mathbf{H}\mathbf{w} = \frac{\eta^2}{2}\mathbf{w}^\top(\mathbf{H} + \lambda\mathbf{I})\mathbf{w},
\end{aligned}$$

where (B.6) holds due to (B.4) and (B.7) holds due to the definition of $\eta$. On the other hand, by the definition of $\mathbf{h}_{\text{opt}}$, $m(\mathbf{h}_{\text{opt}} + \eta\mathbf{w}) - m(\mathbf{h}_{\text{opt}}) \geq 0$. Thus, we have proved (B.2). Finally, we prove (B.3) by showing that

$$\begin{aligned}
\langle \mathbf{g}, \mathbf{h}_{\text{opt}}\rangle + \frac{1}{2}\langle \mathbf{H}\mathbf{h}_{\text{opt}}, \mathbf{h}_{\text{opt}}\rangle + \frac{\theta}{6}\|\mathbf{h}_{\text{opt}}\|_2^3 &= \left\langle \mathbf{g} + \mathbf{H}\mathbf{h}_{\text{opt}} + \frac{\theta}{2}\|\mathbf{h}_{\text{opt}}\|_2\mathbf{h}_{\text{opt}}, \mathbf{h}_{\text{opt}}\right\rangle - \frac{1}{2}\mathbf{h}_{\text{opt}}^\top(\mathbf{H} + \lambda\mathbf{I})\mathbf{h}_{\text{opt}} - \frac{\theta}{12}\|\mathbf{h}_{\text{opt}}\|_2^3 \\
&= -\frac{1}{2}\mathbf{h}_{\text{opt}}^\top(\mathbf{H} + \lambda\mathbf{I})\mathbf{h}_{\text{opt}} - \frac{\theta}{12}\|\mathbf{h}_{\text{opt}}\|_2^3 \tag{B.8} \\
&\leq -\frac{\theta}{12}\|\mathbf{h}_{\text{opt}}\|_2^3, \tag{B.9}
\end{aligned}$$

where (B.8) holds due to (B.1) and (B.9) holds due to (B.2). $\qquad\square$

## B.2  Proof of Lemma A.2

In order to prove Lemma A.2, we need the following useful lemma.

**Lemma B.1.** Suppose $\mathbf{a}_1, \ldots, \mathbf{a}_N$ are i.i.d. and $\mathbb{E}\mathbf{a}_i = 0$, then

$$\mathbb{E}\left\|\frac{1}{N}\sum_{i=1}^N \mathbf{a}_i\right\|_2^{3/2} \leq \frac{1}{N^{3/4}}\left(\mathbb{E}\|\mathbf{a}_i\|_2^2\right)^{3/4}.$$

*Proof of Lemma A.2.* For simplification, we use $\mathbb{E}$ to replace $\mathbb{E}_{\mathbf{v}_{i_t}}$. We have

$$\mathbb{E}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}$$

$$= \mathbb{E}\left\|\frac{1}{b_g}\sum\left[\nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s)\right] + \mathbf{g}^s - \left[\frac{1}{b_g}\sum\nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s) - \mathbf{H}^s\right](\mathbf{x}_t^s - \widehat{\mathbf{x}}^s) - \nabla F(\mathbf{x}_t^s)\right\|_2^{3/2}$$

$$= \mathbb{E}\left\|\frac{1}{b_g}\sum\left[\nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s) - \nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s) - \nabla F(\mathbf{x}_t^s) + \nabla F(\widehat{\mathbf{x}}^s) + \nabla^2 F(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)\right]\right\|_2^{3/2}.$$

Now we set the parameters in Lemma B.1 as

$$N = b_g, \mathbf{a}_{i_t} = \nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s) - \nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s) - \nabla F(\mathbf{x}_t^{s+1}) + \nabla F(\widehat{\mathbf{x}}^s) + \nabla^2 F(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s).$$

We can check that $\mathbf{a}_{i_t}$ satisfy the assumption of Lemma B.1. Thus, by Lemma B.1, we have

$$\mathbb{E}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} \leq \frac{1}{b_g^{3/4}}\Big(\mathbb{E}\big\|\nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s) - \nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)$$

$$- \nabla F(\mathbf{x}_t^s) + \nabla F(\widehat{\mathbf{x}}^s) + \nabla^2 F(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)\big\|_2^2\Big)^{3/4}. \tag{B.10}$$

By Assumption 4.1, we have

$$\left\|\nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s) - \nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s) - \nabla F(\mathbf{x}_t^s) + \nabla F(\widehat{\mathbf{x}}^s) + \nabla^2 F(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)\right\|_2$$

$$\leq \left\|\nabla f_{i_t}(\mathbf{x}_t^s) - \nabla f_{i_t}(\widehat{\mathbf{x}}^s) - \nabla^2 f_{i_t}(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)\right\|_2 + \left\|\nabla F(\mathbf{x}_t^s) - \nabla F(\widehat{\mathbf{x}}^s) - \nabla^2 F(\widehat{\mathbf{x}}^s)(\mathbf{x}_t^s - \widehat{\mathbf{x}}^s)\right\|_2$$

$$\leq \frac{\rho}{2}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2 + \frac{\rho}{2}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2$$

$$= \rho\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2. \tag{B.11}$$

Plugging (B.11) into (B.10) yields

$$\mathbb{E}\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} \leq \frac{1}{b_g^{3/4}}\left(\rho^2\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^4\right)^{3/4} = \frac{\rho^{3/2}}{b_g^{3/4}}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3.$$

$\square$

### B.3  Proof of Lemma A.3

In order to prove Lemma A.3, we need the following supporting lemma.

**Lemma B.2.** Suppose that $q \geq 2, p \geq 2$, and fix $r \geq \max\{q, 2\log p\}$. Consider $\mathbf{Y}_1, ..., \mathbf{Y}_N$ of i.i.d. random self-adjoint matrices with dimension $p \times p$, $\mathbb{E}\mathbf{Y}_i = \mathbf{0}$, then

$$\left[\mathbb{E}\left\|\sum_{i=1}^N \mathbf{Y}_i\right\|_2^q\right]^{1/q} \leq 2\sqrt{er}\left\|\left(\sum_{i=1}^N \mathbb{E}\mathbf{Y}_i^2\right)^{1/2}\right\|_2 + 4er\left(\mathbb{E}\max_i\|\mathbf{Y}_i\|_2^q\right)^{1/q}.$$

*Proof of Lemma A.3.* In the following proof, we replace $\mathbb{E}_{j_t}$ with $\mathbb{E}$ for simplification. We have

$$\mathbb{E}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 = \mathbb{E}\left\|\nabla^2 F(\mathbf{x}_t^s) - \frac{1}{b_h}\left(\sum\left(\nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) + \mathbf{H}^s\right)\right)\right\|_2^3$$

$$= \mathbb{E}\left\|\frac{1}{b_h}\left[\sum\left[\nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) + \mathbf{H}^s - \nabla^2 F(\mathbf{x}_t^s)\right]\right]\right\|_2^3. \tag{B.12}$$

We apply Lemma B.2 with parameters

$$q = 3, p = d, r = 2\log p, \mathbf{Y}_{j_t} = \nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) + \mathbf{H}^s - \nabla^2 F(\mathbf{x}_t^s), N = b_h.$$

It can be easily checked that these parameters satisfy the assumption of Lemma B.2. Meanwhile, by Assumption 4.1, we have the following upper bound for $\mathbf{Y}_{j_t}$:

$$
\begin{aligned}
\left\| \mathbf{Y}_{j_t} \right\|_2 &= \left\| \nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) + \mathbf{H}^s - \nabla^2 F(\mathbf{x}_t^s) \right\|_2 \\
&\leq \left\| \nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) \right\|_2 + \left\| \mathbf{H}^s - \nabla^2 F(\mathbf{x}_t^s) \right\|_2 \\
&\leq \rho \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2 + \rho\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2 \\
&= 2\rho\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2.
\end{aligned}
\tag{B.13}
$$

By Lemma B.2, we have

$$\left[ \mathbb{E} \left\| \sum \mathbf{Y}_{j_t} \right\|_2^3 \right]^{1/3} \leq 2\sqrt{er}\left\| \left( \sum \mathbb{E}\mathbf{Y}_{j_t}^2 \right)^{1/2} \right\|_2 + 4er\left( \mathbb{E} \max_i \|\mathbf{Y}_i\|_2^3 \right)^{1/3}. \tag{B.14}$$

The first term in RHS of (B.14) can be bounded as

$$
\begin{aligned}
2\sqrt{er}\left\| \left( \sum \mathbb{E}\mathbf{Y}_{j_t}^2 \right)^{1/2} \right\|_2 &= 2\sqrt{er}\left\| \sum \mathbb{E}\mathbf{Y}_{j_t}^2 \right\|_2^{1/2} \\
&= 2\sqrt{Ner}\left\| \mathbb{E}\mathbf{Y}_{j_t}^2 \right\|_2^{1/2} \\
&\leq 2\sqrt{Ner}\left( \mathbb{E}\left\| \mathbf{Y}_{j_t}^2 \right\|_2 \right)^{1/2} \\
&= 2\sqrt{Ner}\left( \mathbb{E}\|\mathbf{Y}_{j_t}\|_2^2 \right)^{1/2} \\
&\leq 4\rho\sqrt{Ner}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2,
\end{aligned}
\tag{B.15}
$$

where the first inequality holds due to Jensen's inequality, the third equality holds because $\left\| \mathbf{Y}_{j_t}^2 \right\|_2 = \|\mathbf{Y}_{j_t}\|_2^2$ and the last inequality holds due to (B.13). The second term in RHS of (B.14) can be bounded as

$$4er\left( \mathbb{E} \max_i \left\| \mathbf{Y}_i \right\|_2^3 \right)^{1/3} \leq 4er[(2\rho\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2)^3]^{1/3} = 8\rho er\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2. \tag{B.16}$$

Submitting (B.15), (B.16) into (B.14), we have

$$\left[ \mathbb{E} \left\| \sum \mathbf{Y}_{j_t} \right\|_2^3 \right]^{1/3} \leq 4\rho\sqrt{Ner}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2 + 8\rho er\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2,$$

which immediately implies

$$\mathbb{E}\left\| \frac{1}{N}\sum \mathbf{Y}_{j_t} \right\|_2^3 \leq 64\rho^3\left( \sqrt{\frac{er}{N}} + \frac{2er}{N} \right)^3 \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3. \tag{B.17}$$

Submitting (B.17) into (B.12) with $\mathbf{Y}_{j_t} = \nabla^2 f_{j_t}(\mathbf{x}_t^s) - \nabla^2 f_{j_t}(\widehat{\mathbf{x}}^s) + \mathbf{H}^s - \nabla^2 F(\mathbf{x}_t^s)$, $r = 2\log d$, $N = b_h$, we have

$$
\begin{aligned}
\mathbb{E}\left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^3 &\leq 64\rho^3\left( \sqrt{\frac{2e\log d}{b_h}} + \frac{4e\log d}{b_h} \right)^3 \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \\
&\leq 1200\rho^3\left( \frac{\log d}{b_h} \right)^{3/2} \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3,
\end{aligned}
$$

where the last inequality holds due to $b_h \geq 400\log d$.

$\square$

### B.4   Proof of Lemma A.4

*Proof.* we have

$$
\begin{aligned}
\langle \nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s, \mathbf{h} \rangle &\leq \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 \cdot \|\mathbf{h}\|_2 \\
&= \left( \frac{M^{1/3}}{9^{1/3}} \|\mathbf{h}\|_2 \right) \cdot \left( \frac{9^{1/3}}{M^{1/3}} \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 \right) \\
&\leq \frac{1}{3} \left( \frac{M^{1/3}}{9^{1/3}} \|\mathbf{h}\|_2 \right)^3 + \frac{2}{3} \left( \frac{9^{1/3}}{M^{1/3}} \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 \right)^{3/2} \\
&= \frac{M}{27} \|\mathbf{h}\|_2^3 + \frac{2\|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2}}{M^{1/2}},
\end{aligned}
$$

where the second inequality holds due to Young's inequality. Meanwhile, we have

$$
\begin{aligned}
\langle (\nabla^2 F(\mathbf{x}_t^s) - \mathbf{H}^s)\mathbf{h}, \mathbf{h} \rangle &\leq \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{H}^s \right\|_2 \cdot \|\mathbf{h}\|_2^2 \\
&= \left( \frac{M^{2/3}}{9^{2/3}} \|\mathbf{h}\|_2^2 \right) \cdot \left( \frac{9^{2/3}}{M^{2/3}} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{H}^s \right\|_2 \right) \\
&\leq \frac{2}{3} \left( \frac{M^{2/3}}{9^{2/3}} \|\mathbf{h}\|_2^2 \right)^{3/2} + \frac{1}{3} \left( \frac{9^{2/3}}{M^{2/3}} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{H}^s \right\|_2 \right)^3 \\
&= \frac{2M}{27} \|\mathbf{h}\|_2^3 + \frac{27}{M^2} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^3,
\end{aligned}
$$

where the second inequality holds due to Young's inequality.  $\square$

### B.5   Proof of Lemma A.5

In order to prove Lemma A.5, we need to the following two useful lemmas.

**Lemma B.3.** Under Assumption 4.1, if $M \geq 2\rho$, then we have

$$
\left\| \nabla F(\mathbf{x}_t^s + \mathbf{h}) \right\|_2 \leq M\|\mathbf{h}\|_2^2 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 + \frac{1}{M} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^2 + \left\| \nabla m_t^s(\mathbf{h}) \right\|_2.
$$

**Lemma B.4.** Under Assumption 4.1, if $M \geq 2\rho$, then we have

$$
-\lambda_{\min} \left( \nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) \right) \leq M\|\mathbf{h}\|_2 + \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2 + M \left| \|\mathbf{h}\|_2 - \|\mathbf{h}_t^s\|_2 \right|.
$$

*Proof of Lemma A.5.* By the definition of $\mu$, we can bound $\|\nabla F(\mathbf{x}_t^s + \mathbf{h})\|_2^{3/2}$ and $0 \vee -\rho^{-3/2} \left[ \lambda_{\min} \left( \nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) \right) \right]^3$ separately. To bound $\|\nabla F(\mathbf{x}_t^s + \mathbf{h})\|_2^{3/2}$, applying Lemma B.3 we have

$$
\begin{aligned}
\left\| \nabla F(\mathbf{x}_t^s + \mathbf{h}) \right\|_2^{3/2} &\leq \left[ M\|\mathbf{h}\|_2^2 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 + \frac{1}{M} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^2 + \left\| \nabla m_t^s(\mathbf{h}) \right\|_2 \right]^{3/2} \\
&\leq 2 \left[ M^{3/2}\|\mathbf{h}\|_2^3 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + M^{-3/2} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^3 + \left\| \nabla m_t^s(\mathbf{h}) \right\|_2^{3/2} \right],
\end{aligned}
$$

where the second inequality holds due to the following basic inequality $(a + b + c + d)^{3/2} \leq 2(a^{3/2} + b^{3/2} + c^{3/2} + d^{3/2})$. To bound $-\lambda_{\min} \left( \nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) \right)$, applying Lemma A.3, we have

$$
\begin{aligned}
-\rho^{-3/2} \left[ \lambda_{\min} \left( \nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) \right) \right]^3 &= -C_M^{3/2} M^{-3/2} \left[ \lambda_{\min} \left( \nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) \right) \right]^3 \\
&\leq C_M^{3/2} M^{-3/2} \left[ M\|\mathbf{h}\|_2 + \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2 + M \left| \|\mathbf{h}\|_2 - \|\mathbf{h}_t^s\|_2 \right| \right]^3 \\
&\leq 9 C_M^{3/2} \left[ M^{3/2}\|\mathbf{h}\|_2^3 + M^{-3/2} \left\| \nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s \right\|_2^3 + M^{3/2} \left| \|\mathbf{h}\|_2 - \|\mathbf{h}_t^s\|_2 \right|^3 \right],
\end{aligned}
$$

where the second inequality holds due to $(a + b + c)^3 \leq 9(a^3 + b^3 + c^3)$. Since $9C_M^{3/2} > 2$, we have

$$
\begin{aligned}
&\mu(\mathbf{x}_t^s + \mathbf{h}) \\
&= \max\left\{ \|\nabla F(\mathbf{x}_t^s + \mathbf{h})\|_2^{3/2}, -\rho^{-3/2}\big[\lambda_{\min}\big(\nabla^2 F(\mathbf{x}_t^s + \mathbf{h})\big)\big]^3 \right\} \\
&\leq 9C_M^{3/2}\Big[ M^{3/2}\|\mathbf{h}\|_2^3 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2^{3/2} + M^{-3/2}\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\|_2^3 + \|\nabla m_t^s(\mathbf{h})\|_2^{3/2} + M^{3/2}\big|\|\mathbf{h}\|_2 - \|\mathbf{h}_t^s\|_2\big|^3 \Big],
\end{aligned}
$$

which completes the proof. $\qquad\square$

### B.6 Proof of Lemma A.6

*Proof.* We have

$$
\begin{aligned}
&\|\mathbf{x}_t^s + \mathbf{h} - \widehat{\mathbf{x}}^s\|_2^3 \\
&\leq \big(\|\mathbf{h}\|_2 + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2\big)^3 \\
&= \|\mathbf{h}\|_2^3 + 3\|\mathbf{h}\|_2^2 \cdot \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2 + 3\|\mathbf{h}\|_2 \cdot \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2 + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \\
&= \|\mathbf{h}\|_2^3 + 3\big(T^{1/3}\|\mathbf{h}\|_2^2\big) \cdot \frac{\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2}{T^{1/3}} + 3\big(T^{2/3}\|\mathbf{h}\|_2\big) \cdot \frac{\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2}{T^{2/3}} + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \\
&\leq \|\mathbf{h}\|_2^3 + 3\left(\frac{2}{3}\big(T^{1/3}\|\mathbf{h}\|_2^2\big)^{3/2} + \frac{1}{3}\left(\frac{\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2}{T^{1/3}}\right)^3\right) + 3\left(\frac{1}{3}\big(T^{2/3}\|\mathbf{h}\|_2\big)^3 + \frac{2}{3}\left(\frac{\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^2}{T^{2/3}}\right)^{3/2}\right) + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \\
&= \|\mathbf{h}\|_2^3 + \left(2T^{1/2}\|\mathbf{h}\|_2^3 + \frac{1}{T}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\right) + \left(T^2\|\mathbf{h}\|_2^3 + \frac{2}{T}\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3\right) + \|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|_2^3 \\
&\leq 2T^2\|\mathbf{h}\|_2^3 + \left(1 + \frac{3}{T}\right)\|\mathbf{x}_t^s - \widehat{\mathbf{x}}^s\|^3,
\end{aligned}
\tag{B.18}
$$

where the second inequality holds due to Young's inequality, the last inequality holds because $T \geq 1$. $\qquad\square$

### B.7 Proof of Lemma A.7

*Proof.* By induction, we have for any $0 \leq t \leq T$,

$$
c_t = M\frac{(1 + 3/T)^{T-t} - 1}{1500T^2}.
$$

Then for any $0 \leq t \leq T$,

$$
2c_t T^2 \leq M\frac{2(1 + 3/T)^T}{1500} \leq M\frac{2 \cdot 27}{1500} < \frac{M}{24}.
$$

$\qquad\square$

## C  Proof of Auxiliary Lemmas

In this section, we prove auxillary lemmas used in Appendix B.

### C.1 Proof of Lemma B.1

*Proof.* We have

$$
\mathbb{E}\left\|\frac{1}{N}\sum_{i=1}^N \mathbf{a}_i\right\|_2^{3/2} = \frac{\mathbb{E}\|\sum_{i=1}^N \mathbf{a}_i\|_2^{3/2}}{N^{3/2}} \leq \frac{\big(\mathbb{E}\|\sum_{i=1}^N \mathbf{a}_i\|_2^2\big)^{3/4}}{N^{3/2}} = \frac{\big(\sum_{i=1}^N \mathbb{E}\|\mathbf{a}_i\|_2^2\big)^{3/4}}{N^{3/2}} = \frac{\big(\mathbb{E}\|\mathbf{a}_i\|_2^2\big)^{3/4}}{N^{3/4}}.
\tag{C.1}
$$

The first inequality holds due to Lemma D.1, where we set $s, t$ in Lemma D.1 as $s = 3/2, t = 2$; the second equality holds due to $\mathbb{E}\mathbf{a}_i = 0$ and that $\mathbf{a}_i$ are identically independently distributed. $\qquad\square$

## C.2 Proof of Lemma B.2

*Proof.* This proof is mainly adapted from Chen et al. (2012); Tropp (2016). First, Let $\{\mathbf{Y}_i' : i = 1, \ldots, N\}$ be an independent copy of the sequence $\{\mathbf{Y}_i : i = 1, \ldots, N\}$. We denote $\mathbb{E}_{\mathbf{Y}'}$ to be the expectation over the independent copy $\mathbf{Y}'$. Then $\mathbb{E}_{\mathbf{Y}'}\mathbf{Y}_i' = 0$, then

$$\mathbb{E}\bigg\|\sum_{i=1}^{N}\mathbf{Y}_i\bigg\|_2^q = \mathbb{E}\bigg\|\sum_{i=1}^{N}\mathbb{E}_{\mathbf{Y}'}(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q \leq \mathbb{E}\bigg[\mathbb{E}_{\mathbf{Y}'}\bigg\|\sum_{i=1}^{N}(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q\bigg] = \mathbb{E}\bigg\|\sum_{i=1}^{N}(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q. \tag{C.2}$$

The first equality holds due to $\mathbb{E}_{\mathbf{Y}'}\mathbf{Y}_i' = 0$, the first inequality holds because $\|\cdot\|_2^q$ is a convex function, and the second equality holds because we combine the iterated expectation into a single expectation.

Note that $\mathbf{Y}_i - \mathbf{Y}_i'$ has the same distribution as $\mathbf{Y}_i' - \mathbf{Y}_i$, thus the independent sequence $\{\xi_i(\mathbf{Y}_i - \mathbf{Y}_i') : 1 \leq i \leq n\}$ has the same distribution as $\{\mathbf{Y}_i - \mathbf{Y}_i' : 1 \leq i \leq N\}$, where $\xi_i$ are independent Rademacher random variables, also independent with $\mathbf{Y}_i, \mathbf{Y}_i'$. Therefore,

$$\mathbb{E}\bigg\|\sum_{i=1}^{N}(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q = \mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q. \tag{C.3}$$

Furthermore, we have

$$\mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i(\mathbf{Y}_i - \mathbf{Y}_i')\bigg\|_2^q \leq \mathbb{E}\bigg[2^{q-1}\bigg(\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_2^q + \bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i'\bigg\|_2^q\bigg)\bigg] = 2^q \cdot \mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_2^q. \tag{C.4}$$

The first inequality holds due to $\|\mathbf{A} - \mathbf{B}\|_2^q \leq (\|\mathbf{A}\|_2 + \|\mathbf{B}\|_2)^q \leq 2^{q-1}(\|\mathbf{A}\|_2^q + \|\mathbf{B}\|_2^q)$, where we let $\mathbf{A} = \sum_{i=1}^{N}\xi_i\mathbf{Y}_i, \mathbf{B} = \sum_{i=1}^{N}\xi_i\mathbf{Y}_i'$; the equality holds due to the identical distribution of $\{\xi\mathbf{Y}_i\}$ and $\{\xi\mathbf{Y}_i'\}$. Submitting (C.3), (C.4) into (C.2) yields

$$\mathbb{E}\bigg\|\sum_{i=1}^{N}\mathbf{Y}_i\bigg\|_2^q \leq 2^q \cdot \mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_2^q \tag{C.5}$$

Note that both sides of (C.5) are greater than 0, then we take $q$-th root for both sides, we have

$$\bigg[\mathbb{E}\bigg\|\sum_{i=1}^{N}\mathbf{Y}_i\bigg\|_2^q\bigg]^{1/q} \leq 2\bigg[\mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_2^q\bigg]^{1/q}. \tag{C.6}$$

Next, we have the inequality chain:

$$2\bigg[\mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_2^q\bigg]^{1/q} \leq 2\bigg[\mathbb{E}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_{S_r}^q\bigg]^{1/q}$$
$$= 2\bigg[\mathbb{E}_{\mathbf{Y}_i}\bigg(\mathbb{E}_{\xi_i}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_{S_r}^q\bigg)\bigg]^{1/q}$$
$$\leq 2\bigg[\mathbb{E}_{\mathbf{Y}_i}\bigg(\mathbb{E}_{\xi_i}\bigg\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\bigg\|_{S_r}^r\bigg)^{q/r}\bigg]^{1/q}, \tag{C.7}$$

where the first inequality holds due to $\|\cdot\|_2 \leq \|\cdot\|_{S_r}$, the second inequality holds due to Lyapunov's inequality D.1, where

we set $s, t$ in D.1 as $s = q, t = r$. Since $q < r$, then the second inequality holds. Note we have

$$2\left[\mathbb{E}_{\mathbf{Y}_i}\left(\mathbb{E}_{\xi_i}\left\|\sum_{i=1}^{N}\xi_i\mathbf{Y}_i\right\|_{S_r}^r\right)^{q/r}\right]^{1/q} \leq 2\sqrt{r}\left[\mathbb{E}\left\|\left(\sum_{i=1}^{N}\mathbf{Y}_i^2\right)^{1/2}\right\|_{S_r}^q\right]^{1/q}$$

$$\leq 2\sqrt{r}\left[\mathbb{E}\left(p^{1/r}\left\|\left(\sum_{i=1}^{N}\mathbf{Y}_i^2\right)^{1/2}\right\|_2\right)^q\right]^{1/q}$$

$$\leq 2\sqrt{er}\left[\mathbb{E}\left\|\left(\sum_{i=1}^{N}\mathbf{Y}_i^2\right)^{1/2}\right\|_2^q\right]^{1/q}$$

$$= 2\sqrt{er}\left[\mathbb{E}\left\|\sum_{i=1}^{N}\mathbf{Y}_i^2\right\|_2^{q/2}\right]^{1/q}, \tag{C.8}$$

where the first inequality holds due to Proposition D.2; the second inequality holds because $\|\mathbf{A}\|_{S_r} \leq p^{1/r}\|\mathbf{A}\|_2$, where we set $\mathbf{A} = (\sum_{i=1}^{N}\mathbf{Y}_i^2)^{1/2}$ and $p$ is the dimension of $\mathbf{A}$ ; the third inequality holds because $p^{1/r} \leq p^{1/(2\log p)} = \sqrt{e}$.

Finally, we use Proposition D.3 to bound (C.8). Since $\mathbf{Y}_i^2$ are independent, random, positive-semidefinite matrices, we can set $\mathbf{W}_i$ in Proposition D.3 as $\mathbf{W}_i = \mathbf{Y}_i^2$. Meanwhile, $q/2 \geq 1$, so we have

$$\left[\mathbb{E}\left\|\sum_{i=1}^{N}\mathbf{Y}_i^2\right\|_2^{q/2}\right]^{2/q} \leq \left[\left\|\sum_{i=1}^{N}\mathbb{E}\mathbf{Y}_i^2\right\|_2^{1/2} + 2\sqrt{er}\left(\mathbb{E}\max_i\|\mathbf{Y}_i\|_2^q\right)^{1/q}\right]^2,$$

which immediately implies

$$\left[\mathbb{E}\left\|\sum_{i=1}^{N}\mathbf{Y}_i^2\right\|_2^{q/2}\right]^{1/q} \leq \left\|\sum_{i=1}^{N}\mathbb{E}\mathbf{Y}_i^2\right\|_2^{1/2} + 2\sqrt{er}\left(\mathbb{E}\max_i\|\mathbf{Y}_i\|_2^q\right)^{1/q}. \tag{C.9}$$

Submitting (C.7), (C.8),(C.9) into (C.6), we have the proof completed. $\qquad\square$

## C.3 Proof of Lemma B.3

*Proof.* We have

$$\|\nabla F(\mathbf{x}_t^s + \mathbf{h})\|_2 = \left\|\nabla F(\mathbf{x}_t^s + \mathbf{h}) - \nabla F(\mathbf{x}_t^s) - \nabla^2 F(\mathbf{x}_t^s)\mathbf{h} + \mathbf{v}_t^s + \mathbf{U}_t^s\mathbf{h}\right.$$
$$\left. + \left(\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\right) + \left(\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right)\mathbf{h}\right\|_2$$
$$\leq \left\|\nabla F(\mathbf{x}_t^s + \mathbf{h}) - \nabla F(\mathbf{x}_t^s) - \nabla^2 F(\mathbf{x}_t^s)\mathbf{h}\right\|_2 + \left\|\mathbf{v}_t^s + \mathbf{U}_t^s\mathbf{h} + \frac{M}{2}\|\mathbf{h}\|_2\mathbf{h}\right\|_2$$
$$+ \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 + \left\|\left(\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right)\mathbf{h}\right\|_2 + \frac{M}{2}\|\mathbf{h}\|_2^2, \tag{C.10}$$

where the inequality holds due to triangle inequality. In the following, we are going to bound the right hand side of (C.10). For the first term in the right hand side of (C.10), it can be bounded as

$$\left\|\nabla F(\mathbf{x}_t^s + \mathbf{h}) - \nabla F(\mathbf{x}_t^s) - \nabla^2 F(\mathbf{x}_t^s)\mathbf{h}\right\|_2 \leq \frac{\rho}{2}\|\mathbf{h}\|_2^2 \leq \frac{M}{4}\|\mathbf{h}\|_2^2,$$

where the first inequality holds due to Assumption 4.1 and the second inequality holds due to $2\rho \leq M$. For the second term in the the right hand side of (C.10), it equals to

$$\left\|\mathbf{v}_t^s + \mathbf{U}_t^s\mathbf{h} + \frac{M}{2}\|\mathbf{h}\|_2\mathbf{h}\right\|_2 = \left\|\nabla m_t^s(\mathbf{h})\right\|_2.$$

And the final term can be bounded as

$$\left\|\left(\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right)\mathbf{h}\right\|_2 \leq \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2 \cdot \|\mathbf{h}\|_2 \leq \frac{1}{M}\left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2^2 + \frac{M}{4}\|\mathbf{h}\|_2^2,$$

where the last inequality is due to Young's inequality. Putting all these bounds together and submit them into (C.10), we have

$$\left\|\nabla F(\mathbf{x}_t^s + \mathbf{h})\right\|_2 \leq M\|\mathbf{h}\|_2^2 + \|\nabla F(\mathbf{x}_t^s) - \mathbf{v}_t^s\|_2 + \frac{1}{M}\left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2^2 + \left\|\nabla m_t^s(\mathbf{h})\right\|_2.$$

$\square$

### C.4 Proof of Lemma B.4

*Proof.* We have

$$\begin{aligned}
\nabla^2 F(\mathbf{x}_t^s + \mathbf{h}) &\succeq \nabla^2 F(\mathbf{x}_t^s) - \rho\|\mathbf{h}\|_2\mathbf{I} \\
&\succeq \mathbf{U}_t^s - \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2\mathbf{I} - \rho\|\mathbf{h}\|_2\mathbf{I} \\
&\succeq -\frac{M}{2}\|\mathbf{h}_t^s\|_2\mathbf{I} - \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2\mathbf{I} - \rho\|\mathbf{h}\|_2\mathbf{I},
\end{aligned}$$

where the first inequality holds because $\nabla^2 F$ is $\rho$-Hessian Lipschitz, the last inequality holds due to (A.2) in Lemma A.1. Thus we have

$$\begin{aligned}
-\lambda_{\min}\left(\nabla^2 F(\mathbf{x}_t^s + \mathbf{h})\right) &\leq \frac{M}{2}\|\mathbf{h}_t^s\|_2 + \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2 + \rho\|\mathbf{h}\|_2 \\
&= \frac{M}{2}(\|\mathbf{h}_t^s\|_2 - \|\mathbf{h}\|_2) + \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2 + (\rho + M/2)\|\mathbf{h}\|_2 \\
&\leq M\|\mathbf{h}\|_2 + \left\|\nabla^2 F(\mathbf{x}_t^s) - \mathbf{U}_t^s\right\|_2 + M\big|\|\mathbf{h}_t^s\|_2 - \|\mathbf{h}\|_2\big|,
\end{aligned}$$

where the last inequality holds because $\rho \leq M/2$. $\square$

## D Additional Lemmas and Propositions

**Lemma D.1** (Lyapunov's inequality). (Durrett, 2010) For a random variable $X$, when $0 < s < t$, it holds that

$$(\mathbb{E}|X|^s)^{1/s} \leq (\mathbb{E}|X|^t)^{1/t}.$$

We list two propositions about matrix concentration inequality below. As they play key roles in our next analysis, we use them without proof:

**Proposition D.2** (Matrix Khintchine inequality). (Mackey et al., 2014) Suppose that $r > 2$. Consider a finite sequence $\{\mathbf{A}_i, 1 \leq i \leq N\}$ of deterministic, self-adjoint matrices. Then

$$\left[\mathbb{E}\left\|\sum_{i=1}^N \xi_i \mathbf{A}_i\right\|_{S_r}^r\right]^{1/r} \leq \sqrt{r}\left\|\left[\sum_{i=1}^N \mathbf{A}_i^2\right]^{1/2}\right\|_{S_r},$$

where sequence $\xi_i$ consists of independent Rademacher random variables.

**Proposition D.3.** (Chen et al., 2012) Suppose that $q \geq 1$, and fix $r \geq \max\{q, 2\log p\}$. Consider $\mathbf{W}_1, ..., \mathbf{W}_N$ of independent, random, positive-definite matrices with dimension $p \times p$. Then

$$\left[\mathbb{E}\left\|\sum_{i=1}^N \mathbf{W}_i\right\|_2^q\right]^{1/q} \leq \left[\left\|\sum_{i=1}^N \mathbb{E}\mathbf{W}_i\right\|_2^{1/2} + 2\sqrt{er}\left(\mathbb{E}\max_i \|\mathbf{W}_i\|_2^q\right)^{1/(2q)}\right]^2.$$