

SUPPLEMENT FOR “ACCELERATED GRADIENT-FREE METHOD FOR HEAVILY CONSTRAINED NONCONVEX OPTIMIZATION”

Anonymous authors

Paper under double-blind review

A PROOF OF PROPOSITION 1

Proof 1 Since the $(\mathbf{w}^*, \mathbf{p}^*)$ is the ϵ -stationary point of $\min_{\mathbf{w}} \max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\mathbf{w}, \mathbf{p})$, then we have

$$\|\nabla_{\mathbf{w}} f_0(\mathbf{w}^*) + \beta \sum_{j=1}^m p_j^* 2 \max\{f_j(\mathbf{w}^*), 0\} \nabla_{\mathbf{w}} f_j(\mathbf{w}^*)\|_2^2 \leq \epsilon^2 \quad (1)$$

Let $\alpha_j^* = 2\beta p_j^* \max\{f_j(\mathbf{w}^*), 0\}$ and $\epsilon \leq \epsilon_1$, we have

$$\|\nabla_{\mathbf{w}} f_0(\mathbf{w}^*) + \sum_{j=1}^m \alpha_j^* \nabla_{\mathbf{w}} f_j(\mathbf{w}^*)\|_2^2 \leq \epsilon_1^2 \quad (2)$$

the first condition in Definition 2 is satisfied.

Using $\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}^*, \mathbf{p}^*)\|_2^2 \leq \epsilon$ and $0 \leq p_j^2 \leq 1$, we have

$$\sum_{j=1}^m (\beta \phi_j(\mathbf{w}^*) - \lambda p_j^*)^2 \leq \epsilon^2 \quad (3)$$

Using the inequality $\|a + b\|_2^2 \leq 2\|a\|_2^2 + 2\|b\|_2^2$, we have

$$\begin{aligned} & \frac{1}{2} \beta^2 \sum_{j=1}^m (\max\{f_j(\mathbf{w}^*), 0\})^2 \\ & \leq \beta^2 \sum_{j=1}^m (\max\{f_j(\mathbf{w}^*), 0\} - \lambda p_j^*)^2 + \lambda^2 \sum_{j=1}^m (p_j^*)^2 \\ & \leq \epsilon^2 + m\lambda^2 \end{aligned} \quad (4)$$

Rearrange the above inequality and let $\frac{2\epsilon^2 + 2m\lambda^2}{\beta^2} \leq \epsilon_2^2$, we can obtain

$$\sum_{j=1}^m (\max\{f_j(\mathbf{w}^*), 0\})^2 \leq \epsilon_2^2 \quad (5)$$

Therefore, the second condition in Definition 2 is satisfied.

Based on the inequality $\|\langle \mathbf{a}, \mathbf{b} \rangle\|_2^2 \leq \|\mathbf{a}\|_2^2 \|\mathbf{b}\|_2^2$, we can multiply $\sum_{j=1}^m (\alpha_j^*)^2$ on both sides of the inequality 5, such that we have

$$\left(\sum_{j=1}^m \alpha_j^* \max\{f_j(\mathbf{w}^*), 0\} \right)^2 \leq \sum_{j=1}^m (\alpha_j^*)^2 \sum_{j=1}^m \max\{f_j(\mathbf{w}^*), 0\}^2 \leq \epsilon_2^2 \sum_{j=1}^m (\alpha_j^*)^2 \quad (6)$$

Since $\alpha_j^* \geq 0$ and $\max\{f_j(\mathbf{w}^*), 0\} \geq 0$

$$\sum_{j=1}^m (\alpha_j^* \max\{f_j(\mathbf{w}^*), 0\})^2 \leq \left(\sum_{j=1}^m \alpha_j^* \max\{f_j(\mathbf{w}^*), 0\} \right)^2 \leq \epsilon_2^2 \sum_{j=1}^m (\alpha_j^*)^2 \quad (7)$$

Using inequality 5, we have $(\alpha_j^*)^2 = 4\beta^2(p_j^*)^2(\max\{f_j(\mathbf{w}^*), 0\})^2 \leq 4\beta^2\epsilon_2^2$. Let $4\beta^2\epsilon_2^2 \leq \epsilon_3^2$, we have

$$\sum_{j=1}^m (\alpha_j^* \max\{f_j(\mathbf{w}^*), 0\})^2 \leq \epsilon_3^2 \quad (8)$$

If $f_j(\mathbf{w}^*) \leq 0$, we have $\alpha_j^* = 2\beta p_j^* \max\{f_j(\mathbf{w}^*), 0\} = 0$. Therefore, we have

$$\sum_{j=1}^m (\alpha_j^* f_j(\mathbf{w}^*))^2 \leq \epsilon_3^2, \quad (9)$$

which means that the third condition in Definition 2 is satisfied.

That completes the proof.

B PROOF OF PROPOSITION 2

Proof 2 Assume that a point $\hat{\mathbf{w}}$ satisfies that $\|\nabla_{\mathbf{w}} g(\hat{\mathbf{w}})\|_2 \leq \epsilon$, the optimization problem $\max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{y})$ is strongly concave w.r.t \mathbf{p} and $\mathbf{p}^*(\hat{\mathbf{w}})$ is uniquely defined. Solving this this strongly concave problem $\max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{y})$, we can obtain a point \mathbf{p}' satisfying that

$$\|\nabla_{\mathbf{p}} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{p}')\|_2 \leq \epsilon \text{ and } \|\mathbf{p}' - \mathbf{p}^*(\hat{\mathbf{w}})\|_2 \leq \epsilon. \quad (10)$$

If $\|\nabla_{\mathbf{w}} g(\hat{\mathbf{w}})\|_2 \leq \epsilon$, we have

$$\begin{aligned} & \|\nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{p}')\|_2 \\ & \leq \|\nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{p}') - \nabla_{\mathbf{w}} g(\hat{\mathbf{w}})\|_2 + \|\nabla_{\mathbf{w}} g(\hat{\mathbf{w}})\|_2 \\ & = \|\nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{p}') - \nabla_{\mathbf{w}} \mathcal{L}(\hat{\mathbf{w}}, \mathbf{p}^*(\hat{\mathbf{w}}))\|_2 + \epsilon \\ & \leq L\|\mathbf{p}' - \mathbf{p}^*(\hat{\mathbf{w}})\|_2 + \epsilon \\ & = \mathcal{O}(\epsilon) \end{aligned} \quad (11)$$

C CONVERGENCE ANALYSIS OF DSZOG

Here, we give several lemmas used in our analysis and then give the proof of theorem 1.

Lemma 1 (Nesterov & Spokoiny (2017) Lemma 3) Under Assumption 1, it holds that

$$\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}, \mathbf{p}) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 \leq \frac{\mu^2 L^2(d+3)^3}{4} \quad (12)$$

Lemma 2 (Wang et al. (2020) Lemma 12) Under Assumption 1, for any $\mathbf{w} \in \mathbb{R}^d$, $\mathbf{p}_1, \mathbf{p}_2 \in \Delta^m$, it holds that

$$\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}, \mathbf{p}_1) - \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}, \mathbf{p}_2)\|_2 \leq L^2\|\mathbf{p}_1 - \mathbf{p}_2\|_2^2 \quad (13)$$

Lemma 3 (Nesterov & Spokoiny (2017) Lemma 4) Under Assumption 1, it holds that

$$\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 \leq 2\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}, \mathbf{p})\|_2^2 + \frac{\mu^2 L^2(d+3)^3}{2} \quad (14)$$

Lemma 4 (Lin et al. (2020) Lemma 3.3) The function $g(\cdot) := \max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\cdot, \mathbf{p})$ is $L_g = (L + \kappa L)$ -gradient Lipschitz, and $\nabla g(\mathbf{w}) = \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p}^*(\mathbf{w}))$. Moreover, $\mathbf{p}^*(\mathbf{w}) = \arg \max_{\mathbf{p} \in \Delta^m} \mathcal{L}(\mathbf{w}, \mathbf{p})$ is κ -Lipschitz.

Lemma 5 Under Assumption 1, if $|\mathcal{M}_1| = |\mathcal{M}_2| = M$, it holds that $\mathbb{E}_{\mathbf{u}_{[q]}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}} [\|G_{\mu}^{\mathcal{L}}(\mathbf{w}_t, \mathbf{p}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] \leq 3\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 + \rho(\epsilon, \mu)$, where $\rho(\epsilon, \mu) = \frac{\mu^2 L^2(d+6)^2\epsilon^2}{8} + \frac{\mu^2 L^2(d+3)^3}{2} + \frac{\epsilon^2}{2}$.

Proof 3 Since $\mathbb{E}[G_\mu^\mathcal{L}(\mathbf{w}_t, \mathbf{p}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})] = \nabla_{\mathbf{w}} \mathcal{L}_\mu(\mathbf{w}, \mathbf{p})$ and $|\mathcal{M}_1| = |\mathcal{M}_2| = M$, we have

$$\begin{aligned}
& \mathbb{E}_{\mathbf{u}_{[q]}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}} [\|G_\mu^\mathcal{L}(\mathbf{w}_t, \mathbf{p}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] \\
&= \frac{1}{qM} \mathbb{E}[\|G_\mu^\mathcal{L}(\mathbf{w}_t, \mathbf{p}_t, \ell_i, f_j, \mathbf{u})\|_2^2] \\
&= \frac{1}{qM} \mathbb{E}[\|G_\mu^\mathcal{L}(\mathbf{w}_t, \mathbf{p}_t, \ell_i, f_j, \mathbf{u}) - \nabla_{\mathbf{w}} \mathcal{L}_\mu(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 + \|\nabla_{\mathbf{w}} \mathcal{L}_\mu(\mathbf{w}, \mathbf{p})\|_2^2] \\
&\leq \frac{1}{qM} \mathbb{E}[\|G_\mu^\mathcal{L}(\mathbf{w}_t, \mathbf{p}_t, \ell_i, f_j, \mathbf{u})\|_2^2] + 2\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 + 2\|\nabla_{\mathbf{w}} \mathcal{L}_\mu(\mathbf{w}, \mathbf{p}) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 \\
&\leq \frac{1}{qM} \left[2(d+4) [\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 + \sigma_1^2] + \frac{\mu^2 L^2 (d+6)^3}{2} \right] + 2\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 + \frac{\mu^2 L^2 (d+3)^3}{2} \\
&= \frac{2(d+4)}{qM} \|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 + 2\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 + \frac{2(d+4)\sigma_1^2}{qM} + \frac{\mu^2 L^2 (d+6)^3}{2qM} + \frac{\mu^2 L^2 (d+3)^3}{2}
\end{aligned} \tag{15}$$

The first inequality is due to $\|a+b\|_2^2 \leq 2\|a\|_2^2 + \|b\|_2^2$. The second inequality is due to the Lemma 14 in Nesterov & Spokoiny (2017) and Lemma 1. Let $qM = 4(d+6)(\sigma_1^2 + 1)\epsilon^{-2}$, then we have

$$\mathbb{E}_{\mathbf{u}_{[q]}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}} [\|G_\mu(\mathbf{w}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] \leq 3\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{p})\|_2^2 + \rho(\epsilon, \mu) \tag{16}$$

$$\text{where } \rho(\epsilon, \mu) = \frac{\mu^2 L^2 (d+6)^2 \epsilon^2}{8} + \frac{\mu^2 L^2 (d+3)^3}{2} + \frac{\epsilon^2}{2}.$$

That completes the proof.

Lemma 6 Assume $\{(\mathbf{w}_t, \mathbf{p}_t)\}$ is the sequence generated by our Algorithm 1. By setting $\eta_p = \frac{1}{2L}$, the following inequality holds: $\mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_t\|_2^2] \leq (1 - \frac{1}{4\kappa})\mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1}\|_2^2] + \frac{\sigma^2}{4L^2}$, where $\kappa = \frac{L}{\tau}$.

Proof 4 According to the update rules, we have

$$\begin{aligned}
& \|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_t\|_2^2 \\
&= \|\text{Proj}_{\Delta^m}(\mathbf{p}_{t-1} + \eta_p H(\mathbf{w}_{t-1}, \mathbf{p}_{t-1}) - \mathbf{p}^*(\mathbf{w}_{t-1}))\|_2^2 \\
&\leq \|\mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2 + \eta_p^2 \|H(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})\|_2^2 + 2\eta_p \langle H(\mathbf{w}_{t-1}, \mathbf{p}_{t-1}), \mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1}) \rangle
\end{aligned} \tag{17}$$

For a given t , denote by \mathbb{E} taking expectation with respect to random random samples conditioned on all previous iterations. By taking expectation to both sides of this inequality, we obtain

$$\begin{aligned}
& \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_t\|_2^2] \\
&\leq \mathbb{E}[\|\mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] + \eta_p^2 \mathbb{E}[\|H(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})\|_2^2] - 2\eta_p \langle \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1}), \mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1} \rangle \\
&\leq \mathbb{E}[\|\mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] + \eta_p^2 (\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})\|_2^2 + \sigma_2^2) - 2\eta_p (\mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}^*(\mathbf{w}_{t-1})) - \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})) \\
&\leq \mathbb{E}[\|\mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] + \eta_p^2 (2L(\mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}^*(\mathbf{w}_{t-1})) - \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})) + \sigma_2^2) \\
&\quad - 2\eta_p (\mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}^*(\mathbf{w}_{t-1})) - \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})) \\
&= \mathbb{E}[\|\mathbf{p}_{t-1} - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] - \frac{1}{2L} (\mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}^*(\mathbf{w}_{t-1})) - \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})) + \frac{\sigma_2^2}{4L^2} \\
&\leq \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1}\|_2^2] (1 - \frac{\tau}{4L}) + \frac{\sigma_2^2}{4L^2}
\end{aligned} \tag{18}$$

The second inequality is due to the concavity of \mathcal{L} w.r.t \mathbf{p} . The third inequality is due to the gradient-Lipschitz of \mathcal{L} w.r.t \mathbf{p} and the strongly concave w.r.t \mathbf{p} . The equality is due to $\eta_p = \frac{1}{2L}$. The last is due to strongly concave of \mathcal{L} w.r.t \mathbf{p} .

Lemma 7 Denote $\delta_t = \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2$ and set $\eta_{\mathbf{w}} = \frac{1}{4 \times 16^2 \kappa^2 (\kappa + 1)^2 (L + 1)}$, $\gamma = 1 - \frac{1}{4\kappa} + 3\kappa^2 \eta_{\mathbf{w}}^2 L^2 \leq 1 - \frac{3}{16\kappa} < 1$. It holds that $\mathbb{E}[\delta_t] = \gamma^t \mathbb{E}[\delta_0] + 3\kappa^2 \eta_{\mathbf{w}}^2 \sum_{i=0}^{t-1} \gamma^{t-1-i} \mathbb{E}[\|\nabla g(\mathbf{w}_i)\|_2^2] + \theta_1 \sum_{i=0}^{t-1} \gamma^{t-1-i}$ where $\theta_1 = \kappa^2 \eta_{\mathbf{w}}^2 \rho(\epsilon, \mu) + \frac{\sigma_2^2}{4L^2}$.

Proof 5 We have

$$\begin{aligned}
 \mathbb{E}[\delta_t] &= \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2] \\
 &\leq \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_t\|_2^2] + \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] \\
 &\leq (1 - \frac{1}{4\kappa}) \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1}\|_2^2] + \frac{\sigma_2^2}{4L^2} + \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}^*(\mathbf{w}_{t-1})\|_2^2] \\
 &\leq (1 - \frac{1}{4\kappa}) \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1}\|_2^2] + \frac{\sigma_2^2}{4L^2} + \kappa^2 \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2] \\
 &= \kappa^2 \eta_{\mathbf{w}}^2 \mathbb{E}[\|G_{\mu}^{\mathcal{L}}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] + (1 - \frac{1}{4\kappa}) \mathbb{E}[\delta_{t-1}] + \frac{\sigma_2^2}{4L^2} \quad (19)
 \end{aligned}$$

The third inequality is due to Lemma 4. From Lemma 5, we have

$$\begin{aligned}
 &\mathbb{E}_{\mathbf{u}_{[q]}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}} [\|G_{\mu}(\mathbf{w}_{t-1}, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] \\
 &\leq 3\mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t-1}, \mathbf{p}_{t-1})\|_2^2] + \rho(\epsilon, \mu) \\
 &\leq 3\mathbb{E}[\|\nabla g(\mathbf{w}_{t-1})\|_2^2] + 3l^2 \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_{t-1}) - \mathbf{p}_{t-1}\|_2^2] + \rho(\epsilon, \mu) \quad (20)
 \end{aligned}$$

Then, combining the above inequalities, we have

$$\mathbb{E}[\delta_t] = \gamma^t \mathbb{E}[\delta_0] + 3\kappa^2 \eta_{\mathbf{w}}^2 \sum_{i=0}^{t-1} \gamma^{t-1-i} \mathbb{E}[\|\nabla g(\mathbf{w}_i)\|_2^2] + \theta_1 \sum_{i=0}^{t-1} \gamma^{t-1-i} \quad (21)$$

where $\gamma = 1 - \frac{1}{4\kappa} + 3\kappa^2 \eta_{\mathbf{w}}^2 L^2$ and $\theta_1 = \kappa^2 \eta_{\mathbf{w}}^2 \rho(\epsilon, \mu) + \frac{\sigma_2^2}{4L^2}$.

Then, we give the proof of Theorem 1.

Proof 6

$$\begin{aligned}
& g(\mathbf{w}_{t+1}) \\
\leq & g(\mathbf{w}_t) - \eta_{\mathbf{w}} \langle \nabla g(\mathbf{w}_t), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle + \frac{L_g \eta_{\mathbf{w}}^2}{2} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
= & g(\mathbf{w}_t) - \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)), \\
& + \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t) \\
& + \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle \\
& + \frac{L_g \eta_{\mathbf{w}}^2}{2} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
= & g(\mathbf{w}_t) + \frac{L_g \eta_{\mathbf{w}}^2}{2} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
& + \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle \\
& + \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t) - \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle \\
& - \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle \\
\leq & g(\mathbf{w}_t) + \frac{L_g \eta_{\mathbf{w}}^2}{2} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
& + \frac{1}{L_g} \|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t))\|_2^2 \\
& + \frac{L_g \eta_{\mathbf{w}}^2}{4} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
& + \frac{1}{L_g} \|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t) - \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t))\|_2^2 \\
& + \frac{L_g \eta_{\mathbf{w}}^2}{4} \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 \\
& - \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle \\
\leq & g(\mathbf{w}_t) + L_g \eta_{\mathbf{w}}^2 \|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2 + \frac{\mu^2 L^2 (d+3)^3}{4L_g} + \frac{L^2}{L_g} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
& - \eta_{\mathbf{w}} \langle \nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t), G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]}) \rangle
\end{aligned} \tag{22}$$

where the first inequality is due to Lemma 4. Then second inequality is due to the Young's inequality. The last inequality is due to Lemmas 2 and 1. Then, take expectation w.r.t $\mathbf{u}_{[q]}, \xi_{\mathcal{M}_1}, \zeta_{\mathcal{M}_2}$ on both sides of the above inequality and rearrange it, we obtain

$$\begin{aligned}
& \eta_{\mathbf{w}} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2] \\
\leq & \mathbb{E}[g(\mathbf{w}_t)] - \mathbb{E}[g(\mathbf{w}_{t+1})] + \frac{L^2}{L_g} \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2] + L_g \eta_{\mathbf{w}}^2 \mathbb{E}[\|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] + \frac{\mu^2 L^2 (d+3)^3}{4L_g}
\end{aligned} \tag{23}$$

From Lemma 2, we have

$$\eta_{\mathbf{w}} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}, \mathbf{p}^*(\mathbf{w}))\|_2^2] \leq \eta_{\mathbf{w}} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2] + \eta_{\mathbf{w}} L^2 \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 \tag{24}$$

From Lemma 3, we have

$$\eta_{\mathbf{w}} \|\nabla g(\mathbf{w}_t)\|_2^2 \leq 2\eta_{\mathbf{w}} \|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t))\|_2^2 + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} \tag{25}$$

Combining the above inequalities yields

$$\begin{aligned}
& \eta_{\mathbf{w}} \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \\
& \leq 2\eta_{\mathbf{w}} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}_{\mu}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2] + 2\eta_{\mathbf{w}} L^2 \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} \\
& \leq 2\mathbb{E}[g(\mathbf{w}_t)] - 2\mathbb{E}[g(\mathbf{w}_{t+1})] + \frac{2L_g^2}{L_g} \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2] \\
& \quad + 2L_g \eta_{\mathbf{w}}^2 \mathbb{E}[\|G_{\mu}(\mathbf{w}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})\|_2^2] + \frac{\mu^2 L^2 (d+3)^3}{2L_g} \\
& \quad + 2\eta_{\mathbf{w}} L^2 \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} \\
& \leq 2\mathbb{E}[g(\mathbf{w}_t)] - 2\mathbb{E}[g(\mathbf{w}_{t+1})] + (\frac{2L_g^2}{L_g} + 2\eta_{\mathbf{w}} L^2) \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2] \\
& \quad + \frac{\mu^2 L^2 (d+3)^3}{2L_g} + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} \\
& \quad + 2L_g \eta_{\mathbf{w}}^2 (3\mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] + 3L^2 \mathbb{E}[\|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2] + \rho(\epsilon, \mu)) \\
& = 2\mathbb{E}[g(\mathbf{w}_t)] - 2\mathbb{E}[g(\mathbf{w}_{t+1})] + 6L_g \eta_{\mathbf{w}}^2 \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] + \theta_2 \mathbb{E}[\delta_t] + \theta_3
\end{aligned} \tag{26}$$

where

$$\theta_2 = \frac{2L_g^2}{L_g} + 2\eta_{\mathbf{w}} L^2 + 6L_g \eta_{\mathbf{w}}^2 L^2 \leq 2L + 2\eta_{\mathbf{w}} L^2 + 6\eta_{\mathbf{w}} L^3 (\kappa + 1) \tag{27}$$

and

$$\begin{aligned}
& \theta_3 \\
& = \frac{\mu^2 L^2 (d+3)^3}{2L_g} + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} + 2L_g \eta_{\mathbf{w}}^2 \rho(\epsilon, \mu) \\
& \leq \frac{\mu^2 L (d+3)^3}{2} + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} + 2L(\kappa + 1) \eta_{\mathbf{w}}^2 \rho(\epsilon, \mu) \\
& \leq \frac{\mu^2 L (d+3)^3}{2} + \frac{\eta_{\mathbf{w}} \mu^2 L^2 (d+3)^3}{2} + L^3 (\kappa + 1) \eta_{\mathbf{w}}^2 (\frac{\mu^2 (d+6)^2 \epsilon^2}{4} + \mu^2 (d+3)^3) + L(\kappa + 1) \eta_{\mathbf{w}}^2 \epsilon^2
\end{aligned} \tag{28}$$

Taking sum over $t = 0, \dots, T$, we get

$$\sum_{t=0}^T \mathbb{E}[\delta_t] = \sum_{t=0}^T \gamma^t \mathbb{E}[\delta_0] + 3\kappa^2 \eta_{\mathbf{w}}^2 \sum_{t=0}^T \sum_{i=0}^{t-1} \gamma^{t-1-i} \mathbb{E}[\|\nabla g(\mathbf{w}_i)\|_2^2] + \theta \sum_{t=0}^T \sum_{i=0}^{t-1} \gamma^{t-1-i} \tag{29}$$

Moreover, we have

$$\sum_{t=0}^T \gamma^t \leq 6\kappa \tag{30}$$

$$\sum_{t=0}^T \sum_{i=0}^{t-1} \gamma^{s-1-i} \leq 6\kappa(T+1) \tag{31}$$

$$\sum_{t=0}^T \sum_{i=0}^{t-1} \gamma^{s-1-i} \mathbb{E}[\|\nabla g(\mathbf{w}_i)\|_2^2] \leq 6\kappa \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \tag{32}$$

Thus, we have

$$\sum_{t=0}^T \mathbb{E}[\delta_t] = 6\kappa \mathbb{E}[\delta_0] + 18\kappa^3 \eta_{\mathbf{w}}^2 \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] + \theta_1 6\kappa(T+1) \tag{33}$$

Then, summing 26 over $t = 1, \dots, T$, yields

$$\begin{aligned}
& \eta_w \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \\
&= 2\mathbb{E}[g(\mathbf{w}_0)] - 2\mathbb{E}[g(\mathbf{w}_{T+1})] + 6L_g \eta_w^2 \sum_{t=0}^S \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] + \theta_2 \sum_{t=0}^T \mathbb{E}[\delta_t] + \theta_3(T+1) \\
&\leq 2\mathbb{E}[g(\mathbf{w}_0)] - 2\mathbb{E}[g(\mathbf{w}_{T+1})] + 6L_g \eta_w^2 \sum_{t=0}^S \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \\
&\quad + \theta_2 \left(6\kappa \mathbb{E}[\delta_0] + 18\kappa^3 \eta_w^2 \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] + \theta_1 6\kappa(T+1) \right) \\
&\quad + \theta_3(T+1)
\end{aligned} \tag{34}$$

In addition, we have

$$18\kappa^3 \eta_w^2 \theta_2 \leq \left(\frac{36}{4 \times 16^2} + \frac{36}{4^2 \times 16^4} + \frac{18 \times 6}{4^2 \times 16^4} \right) \eta_w \leq 0.0616 \eta_w \tag{35}$$

and

$$6L_g \eta_w^2 = 6(\kappa + 1)L \eta_w^2 \leq 0.0059 \eta_w \tag{36}$$

Thus, we have

$$0.9325 \eta_w \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \leq 2\mathbb{E}[g(\mathbf{w}_0)] - 2\mathbb{E}[g(\mathbf{w}_{T+1})] + \theta_2 (6\kappa \mathbb{E}[\delta_0] + \theta_1 6\kappa(T+1)) + \theta_3(T+1) \tag{37}$$

Dividing both side of the above inequality by $0.9325(T+1)\eta_w$ and let $\mathbf{p}_0 = \mathbf{p}^*(\mathbf{w}_0)$, we have

$$\frac{1}{T+1} \sum_{t=0}^T \mathbb{E}[\|\nabla g(\mathbf{w}_t)\|_2^2] \leq \frac{2(g(\mathbf{w}_0) - g(\mathbf{w}_{T+1}))}{0.9325(T+1)\eta_w} + \frac{6\theta_2\theta_1\kappa}{0.9325\eta_w} + \frac{\theta_3}{0.9325\eta_w} + \frac{6\theta_2\kappa\sigma_1^2}{0.9325(T+1)\eta_w} \tag{38}$$

Then, we bound the right side by ϵ^2 . Then, we have $T > \max\left\{\frac{2(g(\mathbf{w}_0) - g(\mathbf{w}_{T+1}))}{0.9325\epsilon^2\eta_w}, \frac{\sigma_1^2}{16\kappa\eta_w^2\epsilon^2}\right\}$.

In addititon We have

$$\theta_3 \leq \frac{\mu^2 L(d+3)^3}{2} + \frac{\eta_w \mu^2 L^2(d+3)^3}{2} + L^3(\kappa+1)\eta_w^2 \left(\frac{\mu^2(d+6)^2\epsilon^2}{4} + \mu^2(d+3)^3 \right) + L(\kappa+1)\eta_w^2\epsilon^2 \leq \epsilon^2 \eta_w \tag{39}$$

Let $\mu := \mathcal{O}(\epsilon d^{-3/2} L^{-2})$ and the above inequality holds.

D CONVERGENCE ANALYSIS OF ADSGZO

Here, we give several lemmas used in our analysis and then give the proof of theorem 2.

Lemma 8 Under Assumptions 1 and 3, if $\eta_w L \leq \frac{c_{1,l}}{2c_{1,u}^2}$, we have

$$\begin{aligned}
& g(\mathbf{w}_{t+1}) \\
&\leq g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{4} \|\mathbf{z}_w^t\|_2^2 + \eta_w c_{1,l} \frac{\mu^2 L^2(d+3)^3}{2} + \eta_w c_{1,l} L^2 \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
&\quad + \eta_w c_{1,l} \|\nabla_w \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2
\end{aligned} \tag{40}$$

Proof 7

$$\begin{aligned}
& g(\mathbf{w}_{t+1}) \\
& \leq g(\mathbf{w}_t) + \nabla g(\mathbf{w}_t)^T (\mathbf{w}_{t+1} - \mathbf{w}_t) + \frac{L}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \\
& = g(\mathbf{w}_t) - \eta_w \nabla g(\mathbf{w}_t)^T \frac{\mathbf{z}_w^t}{\sqrt{\|\mathbf{z}_w^t\|_2}} + \frac{L}{2} \frac{\|\mathbf{z}_w^t\|_2^2}{\|\sqrt{\|\mathbf{z}_w^t\|_2}\|_2^2} \\
& \leq g(\mathbf{w}_t) - \eta_w c_{1,l} \nabla g(\mathbf{w}_t)^T \mathbf{z}_w^t + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 \\
& = g(\mathbf{w}_t) + \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t) - \mathbf{z}_w^t\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\mathbf{z}_w^t\|_2^2 + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 \\
& = g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\mathbf{z}_w^t\|_2^2 + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 \\
& \quad + \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t) - \nabla g_\mu(\mathbf{w}_t) + \nabla g_\mu(\mathbf{w}_t) - \nabla \mathbf{L}_\mu(\mathbf{w}_t, \mathbf{p}_t) + \nabla \mathbf{L}_\mu(\mathbf{w}_t, \mathbf{p}_t) - \nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) + \nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2 \\
& \leq g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\mathbf{z}_w^t\|_2^2 + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 \\
& \quad + \eta_w c_{1,l} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla \mathbf{L}_\mu(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t))\|_2^2 + \eta_w c_{1,l} \|\nabla \mathbf{L}_\mu(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) - \nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 \\
& \quad + \eta_w c_{1,l} \|\nabla \mathbf{L}_\mu(\mathbf{w}_t, \mathbf{p}_t) - \nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 + \eta_w c_{1,l} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2 \\
& \leq g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\mathbf{z}_w^t\|_2^2 + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 + \eta_w c_{1,l} \frac{\mu^2 L^2 (d+3)^3}{4} + \eta_w c_{1,l} L^2 \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
& \quad + \eta_w c_{1,l} \frac{\mu^2 L^2 (d+3)^3}{4} + \eta_w c_{1,l} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2 \\
& \leq g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{2} \|\mathbf{z}_w^t\|_2^2 + \frac{L \eta_w^2 c_{1,u}^2}{2} \|\mathbf{z}_w^t\|_2^2 + \eta_w c_{1,l} \frac{\mu^2 L^2 (d+3)^3}{2} + \eta_w c_{1,l} L^2 \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
& \quad + \eta_w c_{1,l} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2 \\
& \leq g(\mathbf{w}_t) - \frac{\eta_w c_{1,l}}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 - \frac{\eta_w c_{1,l}}{4} \|\mathbf{z}_w^t\|_2^2 + \eta_w c_{1,l} \frac{\mu^2 L^2 (d+3)^3}{2} + \eta_w c_{1,l} L^2 \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
& \quad + \eta_w c_{1,l} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_w^t\|_2^2
\end{aligned} \tag{41}$$

The last inequality is due to $\eta_w L \leq \frac{c_{1,l}}{2c_{1,u}^2}$

Lemma 9 Under Assumptions 1 and 3, if $a \leq 1$, $\tau \leq L$ and $\eta_p \leq \frac{1}{3c_{2,l}L}$, we have $\|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \leq -\frac{1}{4a} \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 + \frac{8a\eta_p c_{2,l}}{\tau} \|\nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_p^t\|_2^2 + (1 - \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{8L_g^2}{\tau a \eta_p c_{2,l}} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2$

Proof 8 We have

$$\begin{aligned}
& \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 \\
& \leq \|(1-a)\mathbf{p}_t + a\hat{\mathbf{p}}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 \\
& = \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 + a^2 \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 + 2a \langle \mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t), \mathbf{p}_t - \hat{\mathbf{p}}_{t+1} \rangle
\end{aligned} \tag{42}$$

Thus, we have

$$\langle \mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t), \mathbf{p}_t - \hat{\mathbf{p}}_{t+1} \rangle \geq \frac{1}{2a} (\|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 - \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 - a^2 \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2) \tag{43}$$

Due to the Assumption, we have

$$\mathcal{L}(\mathbf{w}_t, \hat{\mathbf{p}}_{t+1}) \geq \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) + \nabla \mathbf{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t) - \frac{L}{2} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \tag{44}$$

In addition, we have

$$\begin{aligned}
 & \mathcal{L}(\mathbf{w}_t, \mathbf{p}) \\
 & \leq \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) + \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\mathbf{p} - \mathbf{p}_t) - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \\
 & = \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) + \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\mathbf{p} - \hat{\mathbf{p}}_{t+1} + \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t) - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \\
 & = \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) + \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\mathbf{p} - \hat{\mathbf{p}}_{t+1}) + \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t) - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 \\
 & = \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) + (\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t)^T (\mathbf{p} - \hat{\mathbf{p}}_{t+1}) + \langle \mathbf{z}_{\mathbf{p}}^t, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle + \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)^T (\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t) - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2
 \end{aligned} \tag{45}$$

Then, we have

$$\begin{aligned}
 & \mathcal{L}(\mathbf{w}_t, \mathbf{p}) \\
 & \leq \mathcal{L}(\mathbf{w}_t, \hat{\mathbf{p}}_{t+1}) + (\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t)^T (\mathbf{p} - \hat{\mathbf{p}}_{t+1}) + \langle \mathbf{z}_{\mathbf{p}}^t, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 + \frac{L}{2} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2
 \end{aligned} \tag{46}$$

Due to the update rule of $\hat{\mathbf{p}}$, we have

$$\langle \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t - \eta_{\mathbf{p}} \frac{\mathbf{z}_{\mathbf{p}}^t}{\sqrt{\|\mathbf{z}_{\mathbf{p}}^t\|_2}}, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle \geq 0, \quad \forall \mathbf{p} \in \Delta^m \tag{47}$$

Then, we have

$$\begin{aligned}
 & \eta_{\mathbf{p}} c_{2,l} \langle \mathbf{z}_{\mathbf{p}}^t, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle \\
 & \leq \langle \eta_{\mathbf{p}} \frac{\mathbf{z}_{\mathbf{p}}^t}{\sqrt{\|\mathbf{z}_{\mathbf{p}}^t\|_2}}, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle \\
 & \leq \langle \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \hat{\mathbf{p}}_{t+1} \rangle \\
 & = \langle \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \mathbf{p}_t + \mathbf{p}_t - \hat{\mathbf{p}}_{t+1} \rangle \\
 & = - \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \langle \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \mathbf{p}_t \rangle
 \end{aligned} \tag{48}$$

In addition, we have

$$\begin{aligned}
 & \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t, \mathbf{p}^*(\mathbf{w}_t) - \hat{\mathbf{p}}_{t+1} \rangle \\
 & = \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t, \mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t \rangle + \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t, \mathbf{p}_t - \hat{\mathbf{p}}_{t+1} \rangle \\
 & \leq \frac{1}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{1}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 \\
 & \leq \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2
 \end{aligned} \tag{49}$$

Then, we have

$$\begin{aligned}
 & \mathcal{L}(\mathbf{w}_t, \mathbf{p}) \\
 & \leq \mathcal{L}(\mathbf{w}_t, \hat{\mathbf{p}}_{t+1}) - \frac{1}{\eta_{\mathbf{p}} c_{2,l}} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 + \frac{1}{\eta_{\mathbf{p}} c_{2,l}} \langle \hat{\mathbf{p}}_{t+1} - \mathbf{p}_t, \mathbf{p} - \mathbf{p}_t \rangle - \frac{\tau}{2} \|\mathbf{p} - \mathbf{p}_t\|_2^2 + \frac{L}{2} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \\
 & \quad + \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2
 \end{aligned} \tag{50}$$

Let $\mathbf{p} = \mathbf{p}^*(\mathbf{w}_t)$, we have

$$\begin{aligned}
 & \mathcal{L}(\mathbf{w}_t, \hat{\mathbf{p}}_{t+1}) \\
 & \leq \mathcal{L}(\mathbf{w}_t, \mathbf{p}^*(\mathbf{w}_t)) \\
 & \leq \mathcal{L}(\mathbf{w}_t, \hat{\mathbf{p}}_{t+1}) - \frac{1}{\eta_p c_{2,l}} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 - \frac{\tau}{2} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{L}{2} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_t\|_2^2 \\
 & \quad + \frac{2}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{\tau}{4} \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 \\
 & \quad - \frac{1}{2a\eta_p c_{2,l}} (\|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 - \|\mathbf{p}_t - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 - a^2 \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2)
 \end{aligned} \tag{51}$$

Rearrange the inequality, we have

$$\begin{aligned}
 & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 \\
 & \leq -2a\eta_p c_{2,l} \left(\frac{1}{\eta_p c_{2,l}} - \frac{L}{2} - \frac{\tau}{4} - \frac{1}{2b\eta_p c_{2,l}} \right) \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 + \frac{4a\eta_p c_{2,l}}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{2}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 \\
 & \leq -2a\eta_p c_{2,l} \left(\frac{1}{2\eta_p c_{2,l}} - \frac{3L}{4} \right) \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 + \frac{4a\eta_p c_{2,l}}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{2}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2
 \end{aligned} \tag{52}$$

where we use $a \leq 1$, $\tau \leq L$ and $\eta_p \leq \frac{1}{3c_{2,l}L}$.

Then, we have

$$\begin{aligned}
 & \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \\
 & = \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t) + \mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \\
 & \leq (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_t)\|_2^2 + (1 + \frac{4}{\tau a \eta_p c_{2,l}}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \\
 & \leq -2a\eta_p c_{2,l} (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \left(\frac{1}{2\eta_p c_{2,l}} - \frac{3L}{4} \right) \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 + \frac{4a\eta_p c_{2,l}}{\tau} (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{2}) (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + (1 + \frac{4}{\tau a \eta_p c_{2,l}}) L_g^2 \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2 \\
 & \leq -2a\eta_p c_{2,l} (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \left(\frac{1}{2\eta_p c_{2,l}} - \frac{3L}{4} \right) \|\mathbf{p}_t - \hat{\mathbf{p}}_{t+1}\|_2^2 + \frac{8a\eta_p c_{2,l}}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{8L_g^2}{\tau a \eta_p c_{2,l}} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2 \\
 & \leq -\frac{2\eta_p c_{2,l}}{a} (1 + \frac{\tau a \eta_p c_{2,l}}{4}) \left(\frac{1}{2\eta_p c_{2,l}} - \frac{3L}{4} \right) \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 + \frac{8a\eta_p c_{2,l}}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{8L_g^2}{\tau a \eta_p c_{2,l}} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2 \\
 & \leq -\frac{1}{4a} \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 + \frac{8a\eta_p c_{2,l}}{\tau} \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 \\
 & \quad + (1 - \frac{\tau a \eta_p c_{2,l}}{4}) \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \frac{8L_g^2}{\tau a \eta_p c_{2,l}} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2
 \end{aligned} \tag{53}$$

Lemma 10 Under Assumptions 1, 2 and 3, if $b \in (0, 1)$, we have

$$\begin{aligned}
 & \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{p}}^{t+1}\|_2^2] \\
 & \leq (1 - b) \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2] + \frac{1}{b} L^2 [\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b^2 \sigma_2^2
 \end{aligned} \tag{54}$$

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{w}}^{t+1}\|_2^2] \\ & \leq (1-b) \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{w}}^t\|_2^2] + \frac{1}{b} L^2 [\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b^2 \sigma_1^2 \end{aligned} \quad (55)$$

Proof 9

$$\mathbf{z}_{\mathbf{p}}^{t+1} - \mathbf{z}_{\mathbf{p}}^t = -b\mathbf{z}_{\mathbf{p}}^t + bH^{t+1} \quad (56)$$

Then, we have

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{p}}^{t+1}\|_2^2] \\ & = \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{p}}^t - (\mathbf{z}_{\mathbf{p}}^{t+1} - \mathbf{z}_{\mathbf{p}}^t)\|_2^2] \\ & = \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{p}}^t + b\mathbf{z}_{\mathbf{p}}^t - bH^{t+1}\|_2^2] \\ & = \mathbb{E}[\|(1-b)(\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t) + (1-b)(\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)) + b(\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - H^{t+1})\|_2^2] \\ & = (1-b)^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2] + (1-b)^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2] \\ & \quad + b^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - H^{t+1}\|_2^2] + (1-b)^2 \langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t, \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) \rangle \\ & = (1-b)^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2] + (1-b)^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2] \\ & \quad + b^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - H^{t+1}\|_2^2] + (1-b)^2 \mathbb{E}[\langle \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t, \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) \rangle] \\ & \leq (1-b)^2 (1+b) \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2] + (1-b)^2 (1 + \frac{1}{b}) \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t)\|_2^2 \\ & \quad + b^2 \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - H^{t+1}\|_2^2] \\ & \leq (1-b) \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2] + \frac{1}{b} L^2 \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b^2 \sigma_2^2 \end{aligned} \quad (57)$$

Similarly, we have

$$\begin{aligned} & \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{w}}^{t+1}\|_2^2] \\ & \leq (1-b) \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{w}}^t\|_2^2] + \frac{1}{b} L^2 \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b^2 \sigma_1^2 \end{aligned} \quad (58)$$

Then, we give the proof of theorem 2.

Proof 10 Suming up the above inequality, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{w}}^{t+1}\|_2^2] \\ & \leq (1-b) \sum_{t=1}^T \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{w}}^t\|_2^2] + \frac{1}{b} L^2 \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b^2 \sigma_1^2 T \end{aligned} \quad (59)$$

Then, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{w}}^{t+1}\|_2^2] \\ & \leq \frac{1}{b} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{w}}^1\|_2^2] + \frac{1}{b^2} L^2 \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b\sigma_1^2 T \end{aligned} \quad (60)$$

Similarly, we have

$$\begin{aligned} & \sum_{t=1}^T \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_{t+1}, \mathbf{p}_{t+1}) - \mathbf{z}_{\mathbf{p}}^{t+1}\|_2^2] \\ & \leq \frac{1}{b} \mathbb{E}[\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{p}}^1\|_2^2] + \frac{1}{b^2} L^2 \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2] + b\sigma_1^2 T \end{aligned} \quad (61)$$

$$\begin{aligned}
& \sum_{t=1}^T \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \\
& \leq \frac{4}{\tau a \eta_p c_{2,l}} (\|\mathbf{p}^*(\mathbf{w}_1) - \mathbf{p}_1\|_2^2 - \frac{1}{4a} \sum_{t=1}^T \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 + \frac{8a\eta_p c_{2,l}}{\tau} \sum_{t=1}^T \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{8L_g^2}{\tau a \eta_p c_{2,l}} \sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2) \\
& \leq \frac{4}{\tau a \eta_p c_{2,l}} \|\mathbf{p}^*(\mathbf{w}_1) - \mathbf{p}_1\|_2^2 - \frac{1}{\tau a^2 \eta_p c_{2,l}} \sum_{t=1}^T \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 + \frac{32}{\tau^2} \sum_{t=1}^T \|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{p}}^t\|_2^2 + \frac{32L_g^2 \eta_w^2 c_{1,u}^2}{\tau^2 a^2 \eta_p^2 c_{2,l}^2} \sum_{t=1}^T \|\mathbf{z}_{\mathbf{w}}^t\|_2^2
\end{aligned} \tag{62}$$

Thus, we have

$$\begin{aligned}
& \mathbb{E} \left[\sum_{t=1}^T \|\mathbf{p}_{t+1} - \mathbf{p}^*(\mathbf{w}_{t+1})\|_2^2 \right] \\
& \leq \frac{4}{\tau a \eta_p c_{2,l}} \mathbb{E} [\|\mathbf{p}^*(\mathbf{w}_1) - \mathbf{p}_1\|_2^2] + \mathbb{E} \left[\frac{32L_g^2 \eta_w^2 c_{1,u}^2}{\tau^2 a^2 \eta_p^2 c_{2,l}^2} \sum_{t=1}^T \|\mathbf{z}_{\mathbf{w}}^t\|_2^2 - \frac{1}{\tau a^2 \eta_p c_{2,l}} \sum_{t=1}^T \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 \right] \\
& \quad + \frac{32}{\tau^2 b} \mathbb{E} [\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{p}}^1\|_2^2] + \mathbb{E} \left[\sum_{t=1}^T \frac{32}{\tau^2} \left(b\sigma_2^2 + \frac{L^2(\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \|\mathbf{p}_{t+1} - \mathbf{p}_t\|_2^2)}{b^2} \right) \right] \\
& \leq \frac{4}{\tau a \eta_p c_{2,l}} \mathbb{E} [\|\mathbf{p}^*(\mathbf{w}_1) - \mathbf{p}_1\|_2^2] + \mathbb{E} \left[\left(\frac{32L_g^2 \eta_w^2 c_{1,u}^2}{\tau^2 a^2 \eta_p^2 c_{2,l}^2} + \frac{32L^2 \eta_w^2 c_{1,u}^2}{\tau^2 b^2} \right) \sum_{t=1}^T \|\mathbf{z}_{\mathbf{w}}^t\|_2^2 \right] \\
& \quad + \frac{32}{\tau^2 b} \mathbb{E} [\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{p}}^1\|_2^2] + \mathbb{E} \left[\sum_{t=1}^T \frac{32}{\tau^2} b\sigma_2^2 \right] \\
& \quad + \mathbb{E} \left[\left(\frac{32L^2}{\tau^2 b^2} - \frac{1}{\tau a^2 \eta_p c_{2,l}} \right) \sum_{t=1}^T \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 \right]
\end{aligned} \tag{63}$$

In addition, we have

$$\begin{aligned}
& \frac{1}{2} \|\nabla g(\mathbf{w}_t)\|_2^2 \\
& \leq \frac{g(\mathbf{w}_t) - g(\mathbf{w}_{t+1})}{\eta_w c_{1,l}} - \frac{1}{4} \|\mathbf{z}_{\mathbf{w}}^t\|_2^2 + \frac{\mu^2 L^2 (d+3)^3}{2} + L^2 \|\mathbf{p}^*(\mathbf{w}_t) - \mathbf{p}_t\|_2^2 + \|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{p}_t) - \mathbf{z}_{\mathbf{w}}^t\|_2^2
\end{aligned} \tag{64}$$

Summing up from $t = 1, \dots, T$ and taking expectation, we have

$$\begin{aligned}
& \mathbb{E} \left[\frac{1}{2} \sum_{t=1}^T \|\nabla g(\mathbf{w}_t)\|_2^2 \right] \\
& \leq \frac{g(\mathbf{w}_1) - g(\mathbf{w}_{T+1})}{\eta_w c_{1,l}} - \frac{1}{4} \mathbb{E} \left[\sum_{t=1}^T \|\mathbf{z}_{\mathbf{w}}^t\|_2^2 \right] + \frac{\mu^2 T L^2 (d+3)^3}{2} \\
& \quad + \frac{4L^2}{\tau a \eta_p c_{2,l}} \mathbb{E} [\|\mathbf{p}^*(\mathbf{w}_1) - \mathbf{p}_1\|_2^2] + \mathbb{E} \left[\left(\frac{L^2 \eta_w^2 c_{1,u}^2}{b^2} + \frac{32L_g^2 L^2 \eta_w^2 c_{1,u}^2}{\tau^2 a^2 \eta_p^2 c_{2,l}^2} + \frac{32L^4 \eta_w^2 c_{1,u}^2}{\tau^2 b^2} \right) \sum_{t=1}^T \|\mathbf{z}_{\mathbf{w}}^t\|_2^2 \right] \\
& \quad + \frac{32L^2}{\tau^2 b} \mathbb{E} [\|\nabla_{\mathbf{p}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{p}}^1\|_2^2] + \mathbb{E} \left[\sum_{t=1}^T \frac{32L^2}{\tau^2} b\sigma_2^2 \right] \\
& \quad + \mathbb{E} \left[\left(\frac{L^2}{b^2} + \frac{32L^4}{\tau^2 b^2} - \frac{L^2}{\tau a^2 \eta_p c_{2,l}} \right) \sum_{t=1}^T \|\mathbf{p}_t - \mathbf{p}_{t+1}\|_2^2 \right] \\
& \quad + \frac{1}{b} \mathbb{E} [\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_1, \mathbf{p}_1) - \mathbf{z}_{\mathbf{w}}^1\|_2^2] + b\sigma_1^2 T
\end{aligned} \tag{65}$$

Let $\mathbf{p}^*(w_1) = \mathbf{p}_1$, $\mathbf{z}_{\mathbf{p}}^1 = H(\mathbf{w}_t, \mathbf{p}_t, f_{\mathcal{M}_3})$, $\mathbf{z}_{\mathbf{w}}^1 = G_{\mu}^{\mathcal{L}}(\mathbf{w}_t, \mathbf{p}_t, \ell_{\mathcal{M}_1}, f_{\mathcal{M}_2}, \mathbf{u}_{[q]})$, $\eta_p \leq \min\{\frac{b^2}{\tau a^2 c_{2,l}}, \frac{\tau b^2}{32L^2 a^2 c_{2,l}}\}$, $\eta_w^2 \leq \min\{\frac{b^2}{4c_{1,u}^2 L^2}, \frac{\tau^2 a^2 \eta_p^2 c_{2,l}^2}{128L_g^2 L^2 c_{1,u}}, \frac{\tau^2 b^2}{128L^4 c_{1,u}^2}\}$, we have

$$\frac{1}{T} \mathbb{E}[\sum_{t=1}^T \|\nabla g(\mathbf{w}_t)\|_2^2] \leq \frac{2(g(\mathbf{w}_1) - g(\mathbf{w}_{T+1}))}{T \eta_w c_{1,l}} + \frac{64L^2}{\tau^2 b T} \sigma_2^2 + \frac{2\sigma_1^2}{b T} + \mu^2 L^2 (d+3)^3 + \frac{64L^2}{\tau^2} b \sigma_2^2 + 2b\sigma_1^2 \quad (66)$$

Bound the left term by ϵ^2 , we have $\mu \leq \frac{\epsilon}{L(d+3)^{3/2}}$, $b \leq \min\{\frac{\epsilon^2}{2\sigma_1^2}, \frac{\tau^2 \epsilon^2}{64\sigma_2^2 L^2}\}$ and $T \geq \max\{\frac{2(g(\mathbf{w}_1) - g(\mathbf{w}_T))}{\epsilon^2 \eta_w c_{1,l}}, \frac{2\sigma_1^2}{\epsilon^2 b}, \frac{64\sigma_2^2 L^2}{\epsilon^2 \tau^2 b}\}$

REFERENCES

- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
- Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Foundations of Computational Mathematics*, 17(2):527–566, 2017.
- Zhongruo Wang, Krishnakumar Balasubramanian, Shiqian Ma, and Meisam Razaviyayn. Zeroth-order algorithms for nonconvex minimax problems with improved complexities. *stat*, 1050:22, 2020.