

# Appendix for

## *Investigating the Role of Weight Decay in Enhancing Nonconvex SGD*

### Contents

<b>1. Introduction</b>	<b>1</b>
1.1. SGD with Weight Decay . . . . .	1
1.2. Our Contributions . . . . .	2
<b>2. Additional Related Works</b>	<b>2</b>
<b>3. Weight Decay Does not Accelerate SGD</b>	<b>2</b>
3.1. Assumptions for Convergence Analysis . . . . .	2
3.2. Main Convergence Results . . . . .	3
<b>4. Weight Decay Improves Generalization</b>	<b>4</b>
4.1. Setup for Generalization Analysis . . . . .	4
4.2. Main Results for Generalization of SGDW . . . . .	4
<b>5. Extension to Sign-Based Algorithms</b>	<b>5</b>
5.1. Convergence of SignSGDW . . . . .	5
5.2. Generalization of SignSGDW . . . . .	5
<b>6. Numerical Results</b>	<b>6</b>
6.1. Effects of Weight Decay on Convergence of SGD . . . . .	7
6.2. Effects of Weight Decay on Generalization of SGD . . . . .	8
6.3. Numerical Results of SignSGDW . . . . .	8
<b>7. Conclusion</b>	<b>8</b>
<b>A Technical Lemmas</b>	<b>11</b>
<b>B Proof of Theorem 1</b>	<b>12</b>
<b>C Proof of Theorem 2</b>	<b>14</b>
<b>D Proof of Theorem 3</b>	<b>15</b>
<b>E Proof of Theorem 4</b>	<b>16</b>
<b>F. Proofs of Technical Lemmas</b>	<b>17</b>
F.1. Proof of Lemma 2 . . . . .	17
F.2. Proof of Lemma 3 . . . . .	18
F.3. Proof of Lemma 4 . . . . .	18
F.4. Proof of Lemma 5 . . . . .	19
F.5. Proof of Lemma 6 . . . . .	19

### A. Technical Lemmas

**Lemma 3** Assume  $(\mathbf{w}^t)_{t \geq 0}$  be generated by SGDW, as  $0 < \alpha < 1$ , we have

$$\sum_{t=1}^T \mathbb{E} \|\mathbf{w}^t\|^2 \leq \frac{\gamma^2(1 - \alpha^T)}{(1 - \alpha)^2} \sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t)\|^2 + \frac{\gamma^2 T(1 - \alpha^T)}{1 - \alpha} \sigma^2,$$

**Lemma 4** Let  $(\mathbf{w}^t)_{t \geq 0}$  be generated by SignSGDW, it holds that

$$\|\mathbf{w}^t\| \leq \frac{1-\alpha^T}{1-\alpha} \gamma \sqrt{d} \leq \frac{1}{1-\alpha} \gamma \sqrt{d}, \quad \|\mathbf{w}^t\|_\infty \leq \frac{1-\alpha^T}{1-\alpha} \gamma,$$

provided  $0 \leq t \leq T$ .

**Lemma 5** Let  $\mathbf{x}^\dagger, \mathbf{m} \in \mathbb{R}^d$  be arbitrary vectors such that  $\|\mathbf{x}^\dagger\| \leq \frac{1}{1-\alpha} \gamma \sqrt{d}$ ,  $\|\mathbf{x}^\dagger\|_\infty \leq \frac{1-\alpha^T}{1-\alpha} \gamma$ , and

$$\mathbf{x}^\dagger = \alpha \mathbf{x}^\dagger - \gamma \text{Sign}(\mathbf{m}), \quad (3)$$

and  $\boldsymbol{\epsilon} := \mathbf{m} - \nabla f_S(\mathbf{x}^\dagger)$ . If  $|\langle \nabla f_S(\mathbf{x}^\dagger), \mathbf{x}^\dagger \rangle| \leq \delta \|\nabla f_S(\mathbf{x}^\dagger)\|_1 \cdot \|\mathbf{x}^\dagger\|_\infty$ , then we have

$$f_S(\mathbf{x}^\dagger) - f_S(\mathbf{x}^\dagger) \leq -\gamma [1 - (1-\alpha^T)\delta] \|\nabla f_S(\mathbf{x}^\dagger)\|_1 + 2\sqrt{d}\gamma\|\boldsymbol{\epsilon}\| + 2Ld\gamma^2.$$

**Lemma 6** Let Assumption 1, 2, and 3 hold, and let  $(\mathbf{x}^t)_{t \geq 1}$  and  $(\mathbf{y}^k)_{t \geq 1}$  be the output of SignSGDW on these two neighbor sets of  $n$  samples with the same initial points. Then:

1) In the general case, it holds that

$$\sup_{\xi \in \Xi} \sqrt{\mathbb{E}_{\mathcal{A}_{\text{SignSGDW}}} \|\nabla f(\mathcal{A}_{\text{SignSGDW}}(S); \xi) - \nabla f(\mathcal{A}_{\text{SignSGDW}}(S'); \xi)\|^2} = \mathcal{O}\left(\frac{T^{3/4}}{\sqrt{n}}\right).$$

2) If  $\delta < 1$  and  $\gamma = \frac{1}{L\sqrt{d}T^{3/4}}$ ,  $\alpha = 1 - \frac{1}{T^{1/4}}$ , then we have

$$\sup_{\xi \in \Xi} \sqrt{\mathbb{E}_{\mathcal{A}_{\text{SignSGDW}}} \|\nabla f(\mathcal{A}_{\text{SignSGDW}}(S); \xi) - \nabla f(\mathcal{A}_{\text{SignSGDW}}(S'); \xi)\|^2} = \mathcal{O}\left(\frac{L\gamma\sqrt{d}T^{3/4}}{\sqrt{n}}\right) = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right).$$

## B. Proof of Theorem 1

The Lipschitz property gives us

$$\begin{aligned} f_S(\mathbf{w}^{t+1}) &\leq f_S(\mathbf{w}^t) + \langle \nabla f_S(\mathbf{w}^t), \mathbf{w}^{t+1} - \mathbf{w}^t \rangle + \frac{L^2}{2} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \\ &\leq f_S(\mathbf{w}^t) + \langle \nabla f_S(\mathbf{w}^t), \mathbf{w}^{t+1} - \alpha \mathbf{w}^t + (\alpha - 1) \mathbf{w}^t \rangle + \frac{L^2}{2} \|\mathbf{w}^{t+1} - \alpha \mathbf{w}^t + (\alpha - 1) \mathbf{w}^t\|^2 \\ &= f_S(\mathbf{w}^t) + \langle \nabla f_S(\mathbf{w}^t), (\alpha - 1) \mathbf{w}^t - \gamma \mathbf{g}^t \rangle + \frac{L^2}{2} \|(\alpha - 1) \mathbf{w}^t - \gamma \mathbf{g}^t\|^2 \\ &= f_S(\mathbf{w}^t) - \gamma \langle \mathbf{g}^t, \nabla f_S(\mathbf{w}^t) \rangle - (1 - \alpha) \langle \nabla f_S(\mathbf{w}^t), \mathbf{w}^t \rangle + \frac{L^2 \gamma^2}{2} \|\mathbf{g}^t\|^2 \\ &\quad + (1 - \alpha) \gamma L^2 \langle \mathbf{g}^t, \mathbf{w}^t \rangle + \frac{(1 - \alpha)^2 L^2}{2} \|\mathbf{w}^t\|^2. \end{aligned}$$

Taking the expectation of both sides over  $\chi^t$ , we get

$$\begin{aligned} \mathbb{E}[f_S(\mathbf{w}^{t+1}) | \chi^t] &\leq f_S(\mathbf{w}^t) - \gamma \|\nabla f_S(\mathbf{w}^t)\|^2 - (1 - \alpha)(1 - \gamma L^2) \langle \nabla f_S(\mathbf{w}^t), \mathbf{w}^t \rangle \\ &\quad + \frac{(1 - \alpha)^2 L^2}{2} \|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2}{2} \mathbb{E}[\|\mathbf{g}^t\|^2 | \chi^t] \\ &\leq f_S(\mathbf{w}^t) - \gamma \left(1 - \frac{L\gamma}{2}\right) \|\nabla f_S(\mathbf{w}^t)\|^2 + (1 - \alpha)(1 - \gamma L^2) \underbrace{\langle -\nabla f_S(\mathbf{w}^t), \mathbf{w}^t \rangle}_{(\dagger)} \\ &\quad + \frac{(1 - \alpha)^2 L^2}{2} \|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2 \sigma^2}{2}. \end{aligned} \quad (4)$$

**I.** Using  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \|\mathbf{a}\| \|\mathbf{b}\| \leq \frac{\|\mathbf{a}\|^2}{2\sqrt{1-\alpha}} + \frac{\sqrt{1-\alpha}\|\mathbf{b}\|^2}{2}$  for ( $\dagger$ ) in (4) as  $\nabla f_S(\mathbf{w}^t) \rightarrow \mathbf{a}$  and  $\mathbf{w}^t \rightarrow \mathbf{b}$ , we get

$$\begin{aligned} \mathbb{E}[f_S(\mathbf{w}^{t+1}) | \chi^t] &\leq f_S(\mathbf{w}^t) - \left[ \gamma(1 - \frac{L\gamma}{2}) - \frac{\sqrt{1-\alpha}(1-\gamma L^2)}{2} \right] \|\nabla f_S(\mathbf{w}^t)\|^2 \\ &\quad + \left[ \frac{(1-\alpha)^2 L^2}{2} + \frac{(1-\alpha)^{3/2}(1-\gamma L^2)}{2} \right] \|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2 \sigma^2}{2}, \end{aligned} \tag{5}$$

We use some shorthand notation for constants as follows

$$C_1 := \gamma(1 - \frac{L\gamma}{2}) - \frac{\sqrt{1-\alpha}(1-\gamma L^2)}{2}, \quad C_2 := \frac{(1-\alpha)^2 L^2}{2} + \frac{(1-\alpha)^{3/2}(1-\gamma L^2)}{2}.$$

Taking full expectations on both sides of (5) results in

$$\mathbb{E}f_S(\mathbf{w}^{t+1}) \leq \mathbb{E}f_S(\mathbf{w}^t) - C_1 \mathbb{E}\|\nabla f_S(\mathbf{w}^t)\|^2 + C_2 \mathbb{E}\|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2 \sigma^2}{2}.$$

Summing the inequality from  $t = 1$  to  $T$  and using Lemma 3, we get

$$\left[ C_1 - C_2 \frac{\gamma^2(1-\alpha^T)}{(1-\alpha)^2} \right] \sum_{t=1}^T \|\nabla f_S(\mathbf{w}^t)\|^2 \leq f_S(\mathbf{w}^1) - \min f_S + C_2 \frac{\gamma^2 \sigma^2 T}{1-\alpha} + \frac{L^2 \gamma^2 \sigma^2 T}{2}.$$

Noticing that  $\gamma$  is very small and  $1 - \gamma^2 \leq \alpha < 1$ , it follows

$$C_1 \geq \frac{\gamma}{2}, \quad C_2 \approx \frac{(1-\alpha)^{3/2}}{2}, \quad \frac{1-\alpha^T}{1-\alpha} \approx T.$$

and

$$C_2 \frac{\gamma^2(1-\alpha^T)}{2(1-\alpha)^2} \approx \frac{\gamma^2(1-\alpha^T)}{2\sqrt{1-\alpha}} = \frac{\gamma^2(1-\alpha^T)\sqrt{1-\alpha}}{2(1-\alpha)} \approx \frac{\gamma^2 T \sqrt{1-\alpha}}{2} \leq \gamma/4.$$

Based on these approximations, as  $\gamma$  is small, we can see that

$$C_1 - C_2 \frac{\gamma^2(1-\alpha^T)}{(1-\alpha)^2} \geq \frac{\gamma}{4}.$$

Therefore, we obtain the following bound for the rate

$$\frac{\sum_{t=1}^T \mathbb{E}\|\nabla f_S(\mathbf{w}^t)\|^2}{T} = \mathcal{O}\left(\frac{\mathbb{E}f_S(\mathbf{w}^1) - \min f_S}{\gamma T} + \frac{L^2 \gamma \sigma^2}{2} + \sqrt{1-\alpha} \gamma \sigma^2\right).$$

**II.** If  $0 \leq \delta < 1$ : Using  $\langle \mathbf{a}, \mathbf{b} \rangle \leq \delta_T \|\mathbf{a}\| \|\mathbf{b}\| \leq \frac{\delta_T \|\mathbf{a}\|^2}{2} + \frac{\delta_T \|\mathbf{b}\|^2}{2}$  for ( $\dagger$ ) in (4) as  $\nabla f_S(\mathbf{w}^t) \rightarrow \mathbf{a}$  and  $\mathbf{w}^t \rightarrow \mathbf{b}$ , and using the fact that  $\langle -\nabla f_S(\mathbf{w}^t), \mathbf{w}^t \rangle \leq \|\nabla f_S(\mathbf{w}^t)\| \cdot \|\mathbf{w}^t\|$ , we get

$$\begin{aligned} \mathbb{E}[f_S(\mathbf{w}^{t+1}) | \chi^t] &\leq f_S(\mathbf{w}^t) - \left[ \gamma(1 - \frac{L\gamma}{2}) - \frac{\delta_T(1-\alpha)(1-\gamma L^2)}{2} \right] \|\nabla f_S(\mathbf{w}^t)\|^2 \\ &\quad + \left[ \frac{(1-\alpha)^2 L^2}{2} + \frac{\delta_T(1-\alpha)(1-\gamma L^2)}{2} \right] \|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2 \sigma^2}{2}, \end{aligned} \tag{6}$$

We use some shorthand notation for constants as follows

$$C_1 := \gamma(1 - \frac{L\gamma}{2}) - \frac{\delta_T(1-\alpha)(1-\gamma L^2)}{2}, \quad C_2 := \frac{(1-\alpha)^2 L^2}{2} + \frac{\delta_T(1-\alpha)(1-\gamma L^2)}{2}.$$

Taking full expectations on both sides of (4), we get

$$\mathbb{E}f_S(\mathbf{w}^{t+1}) \leq \mathbb{E}f_S(\mathbf{w}^t) - C_1 \mathbb{E}\|\nabla f_S(\mathbf{w}^t)\|^2 + C_2 \mathbb{E}\|\mathbf{w}^t\|^2 + \frac{L^2 \gamma^2 \sigma^2}{2}.$$

Summing the inequality from  $t = 1$  to  $T$  and using Lemma 3, we get

$$\left[ C_1 - C_2 \frac{\gamma^2(1 - \alpha^T)}{(1 - \alpha)^2} \right] \sum_{t=1}^T \|\nabla f_S(\mathbf{w}^t)\|^2 \leq f_S(\mathbf{w}^1) - \min f_S + C_2 \frac{\gamma^2(1 - \alpha^T)\sigma^2 T}{1 - \alpha} + \frac{L^2\gamma^2\sigma^2 T}{2}.$$

Noticing that  $\gamma$  is very small and  $1 - \alpha = \gamma$ , it follows

$$C_1 \approx \gamma - \frac{\delta_T(1 - \alpha)}{2} \geq \gamma - \frac{\delta_T\gamma}{2}, \quad C_2 \approx \frac{\delta_T(1 - \alpha)}{2} \leq \frac{\delta_T\gamma}{2},$$

and

$$C_2 \frac{\gamma^2(1 - \alpha^T)}{(1 - \alpha)^2} \approx \delta_T \frac{\gamma^2(1 - \alpha^T)}{2(1 - \alpha)} \leq \frac{\delta_T\gamma^2}{2(1 - \alpha)} = \frac{\delta_T\gamma}{2}.$$

Therefore, we are led to

$$C_1 - C_2 \frac{\gamma^2(1 - \alpha^T)}{(1 - \alpha)^2} \geq (1 - \delta_T)\gamma.$$

Hence, we get

$$\frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|^2}{T} = \mathcal{O}\left(\frac{\mathbb{E} f_S(\mathbf{w}^1) - \min f_S}{(1 - \delta_T)\gamma T} + \frac{(L^2 + 1)\gamma\sigma^2}{2(1 - \delta_T)}\right).$$

Combining the two cases above, we then complete the proof.

## C. Proof of Theorem 2

With the Cauchy's inequality, it holds that

$$\begin{aligned} \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|}{T} &= \frac{\sum_{t=1}^T \mathbb{E}(1 \cdot \|\nabla f_S(\mathbf{w}^t)\|)}{T} \\ &\leq \frac{\sqrt{T \sum_{i=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|^2}}{T} = \frac{\sqrt{\sum_{i=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|^2}}{\sqrt{T}}. \end{aligned}$$

I. As  $1 - \alpha \leq \frac{1}{2T}$ , from Theorem 1, it holds that  $\sum_{i=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|^2 = \mathcal{O}(\sqrt{T})$ , indicating the following result

$$\frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|}{T} = \mathcal{O}\left(\frac{1}{T^{1/4}}\right).$$

From Lemmas 2 and 1, we can see that

$$\mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\| = \mathcal{O}\left(\sqrt{\frac{t}{n}}\right) = \mathcal{O}\left(\sqrt{\frac{T}{n}}\right).$$

Based on the inequality above, we have

$$\frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} = \mathcal{O}\left(\sqrt{\frac{T}{n}}\right).$$

By the triangle inequality, we have

$$\begin{aligned} \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t)\|}{T} &= \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|}{T} + \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} \\ &= \mathcal{O}\left(\frac{1}{T^{1/4}} + \sqrt{\frac{T}{n}}\right). \end{aligned}$$

**II.** For the case  $1 - \alpha = \frac{1}{\sqrt{T}}$ , from Theorem 1, we can get

$$\frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|}{T} = \mathcal{O}\left(\frac{1}{\sqrt{1 - \delta_T} T^{1/4}}\right).$$

In this case, we also have

$$\frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} = \mathcal{O}\left(\sqrt{\frac{\gamma T}{n}}\right) = \mathcal{O}\left(\frac{T^{1/4}}{\sqrt{n}}\right).$$

Therefore, we can get the following generalization bound

$$\begin{aligned} \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t)\|}{T} &= \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|}{T} + \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} \\ &= \mathcal{O}\left(\frac{1}{\sqrt{1 - \delta_T} T^{1/4}} + \frac{T^{1/4}}{\sqrt{n}}\right). \end{aligned}$$

Combination of the two cases above then completes the proof.

## D. Proof of Theorem 3

Let

$$\mathbf{g}^t := \nabla f(\mathbf{w}^t; \xi_{i_t}), \boldsymbol{\epsilon}^t := \mathbf{m}^t - \nabla f_S(\mathbf{w}^t), \boldsymbol{\delta}^t := \mathbf{g}^t - \nabla f_S(\mathbf{w}^t).$$

Through direct computations, we obtain

$$\begin{aligned} \mathbf{m}^t &= \theta \mathbf{m}^{t-1} + (1 - \theta) \mathbf{g}^t = \theta(\boldsymbol{\epsilon}^{t-1} + \nabla f_S(\mathbf{w}^{t-1})) + (1 - \theta)(\boldsymbol{\delta}^t + \nabla f_S(\mathbf{w}^t)) \\ \Rightarrow \boldsymbol{\epsilon}^t &= \mathbf{m}^t - \nabla f_S(\mathbf{w}^t) = \theta \boldsymbol{\epsilon}^{t-1} + \theta(\nabla f_S(\mathbf{w}^{t-1}) - \nabla f_S(\mathbf{w}^t)) + (1 - \theta) \boldsymbol{\delta}^t. \end{aligned}$$

Noting the smoothness of the gradient, we have

$$\begin{aligned} \|\boldsymbol{s}^t\| &= \|\nabla f_S(\mathbf{w}^{t-1}) - \nabla f_S(\mathbf{w}^t)\| \leq L \|\mathbf{w}^{t-1} - \mathbf{w}^t\| \\ &= L \|\alpha \mathbf{w}^{t-1} - \mathbf{w}^t + (1 - \alpha) \mathbf{w}^{t-1}\| \leq L \|\alpha \mathbf{w}^{t-1} - \mathbf{w}^t\| + L \|(\alpha - 1) \mathbf{w}^{t-1}\| \\ &\leq L \sqrt{d} \gamma + (1 - \alpha) L \|\mathbf{w}^{t-1}\| \leq 2L \sqrt{d} \gamma. \end{aligned}$$

This leads us to

$$\boldsymbol{\epsilon}^t = \theta \boldsymbol{\epsilon}^{t-1} + \boldsymbol{s}^t + (1 - \theta) \boldsymbol{\delta}^t.$$

Using Mathematical Induction, we derive the following result

$$\boldsymbol{\epsilon}^t = \theta^t \boldsymbol{\epsilon}^0 + \theta \sum_{i=1}^t \theta^{t-i} \boldsymbol{s}^i + (1 - \theta) \sum_{i=1}^t \theta^{t-i} \boldsymbol{\delta}^i. \quad (7)$$

Taking norms on both sides gives

$$\|\boldsymbol{\epsilon}^t\| \leq 2L \sqrt{d} \gamma \sum_{i=1}^t \theta^i + (1 - \theta) \left\| \sum_{i=1}^t \theta^{t-i} \boldsymbol{\delta}^i \right\| + \theta^t \|\boldsymbol{\epsilon}^0\|.$$

Furthermore, by taking expectations on both sides, we have

$$\mathbb{E} \|\boldsymbol{\epsilon}^t\| \leq \frac{2L \sqrt{d} \gamma \theta}{1 - \theta} + (1 - \theta) \mathbb{E} \left\| \sum_{i=1}^t \theta^{t-i} \boldsymbol{\delta}^i \right\| + \theta^t \|\boldsymbol{\epsilon}^0\|. \quad (8)$$

Applying Cauchy's inequality and considering the independence of  $(s^i)_{1 \leq i \leq t}$ , together with the fact that  $\mathbb{E}|\delta^i|^2 \leq \sigma^2$ , we have

$$\mathbb{E} \left\| \sum_{i=1}^t \theta^{t-i} \delta^i \right\| \leq \sqrt{\mathbb{E} \left\| \sum_{i=1}^t \theta^{t-i} \delta^i \right\|^2} = \sqrt{\mathbb{E} \sum_{i=1}^t \theta^{2t-2i} \|\delta^i\|^2} \leq \frac{\sigma}{\sqrt{1-\theta^2}}.$$

Consequently, we obtain

$$\mathbb{E} \|\epsilon^t\| \leq \frac{2L\sqrt{d}\gamma\theta}{1-\theta} + \frac{\sqrt{1-\theta}}{\sqrt{1+\theta}}\sigma + \theta^t \|\epsilon^0\| \leq \frac{2L\sqrt{d}\gamma}{1-\theta} + \sqrt{1-\theta}\sigma + \theta^t \|\epsilon^0\|.$$

Leveraging Lemma 5 with  $\mathbf{x}^\dagger \rightarrow \mathbf{x}^t$  and  $\mathbf{m} \rightarrow \mathbf{m}^t$  and Lemma 4, we derive

$$f_S(\mathbf{w}^{t+1}) - f_S(\mathbf{w}^t) \leq -\gamma[1 - (1 - \alpha^T)\delta_T] \|\nabla f_S(\mathbf{w}^t)\|_1 + 2\sqrt{d}\gamma \|\epsilon^t\| + 2Ld\gamma^2.$$

Summing the recursion from  $t = 1$  to  $T$  yields

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{x}^t)\|_1 &\leq \frac{f_S(\mathbf{w}^1) - \min f_S}{\gamma[1 - (1 - \alpha^T)\delta_T]T} + \frac{4Ld\gamma}{[1 - (1 - \alpha^T)\delta_T](1-\theta)} + \frac{2\sqrt{d}\sqrt{1-\theta}\sigma}{[1 - (1 - \alpha^T)\delta_T]} \\ &\quad + \frac{2L\gamma d}{[1 - (1 - \alpha^T)\delta_T]} + \frac{2\sqrt{d}\|\epsilon^0\|}{[1 - (1 - \alpha^T)\delta_T](1-\theta)T}. \end{aligned}$$

Let  $1 - \theta = \frac{1}{\sqrt{T}}$  and  $\gamma = \frac{1}{L\sqrt{d}T^{3/4}}$ , we arrive at

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{x}^t)\|_1 = \mathcal{O} \left( \frac{\sqrt{d}(f_S(\mathbf{w}^1) - \min f_S + \sigma + \|\nabla f_S(\mathbf{w}^0)\|)}{[1 - (1 - \alpha^T)\delta_T]T^{1/4}} \right).$$

Recalling Lemma 4, we observe that the weight decay is approximately  $\frac{1}{T^{3/4}}$  when  $\gamma = \frac{1}{L\sqrt{d}T^{3/4}}$ . Otherwise, the optimization will be confined to a very small region around  $\mathbf{0}$ .

**I.** As  $1 - \frac{1}{T} \leq \alpha < 1$ ,  $\alpha^T \asymp 1$  yielding

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{x}^t)\|_1 = \mathcal{O} \left( \frac{\sqrt{d}(f_S(\mathbf{w}^1) - \min f_S + \sigma + \|\nabla f_S(\mathbf{w}^0)\|)}{T^{1/4}} \right).$$

**II.** As  $1 - \frac{1}{T^{3/4}} = \alpha$ ,  $\alpha^T \asymp 0$ , we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{x}^t)\|_1 &= \mathcal{O} \left( \frac{\sqrt{d}(f_S(\mathbf{w}^1) - \min f_S + \sigma + \|\nabla f_S(\mathbf{w}^0)\|)}{[(1 - \delta_T) + \delta_T\alpha^T]T^{1/4}} \right) \\ &= \mathcal{O} \left( \frac{\sqrt{d}(f_S(\mathbf{w}^1) - \min f_S + \sigma + \|\nabla f_S(\mathbf{w}^0)\|)}{(1 - \delta_T)T^{1/4}} \right). \end{aligned}$$

Note that  $\mathbf{w}^1 = \mathbf{w}^0$  because  $\mathbf{m}^0 = \mathbf{0}$ , we then proved the result.

## E. Proof of Theorem 4

**I.** Using Lemmas 1 and 6, in the general case, we then get

$$\begin{aligned} \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t)\|_1}{T} &\leq \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|_1}{T} + \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|_1}{T} \\ &\leq \frac{\sum_{t=1}^T \sqrt{d} \mathbb{E} \|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} + \frac{\sum_{t=1}^T \mathbb{E} \|\nabla f_S(\mathbf{w}^t)\|_1}{T} \\ &= \mathcal{O} \left( \frac{\sqrt{d}}{T^{1/4}} + \frac{\sqrt{dT^{3/4}}}{\sqrt{n}} \right). \end{aligned}$$

**II.** If  $\delta < 1$ , we have

$$\begin{aligned} \frac{\sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}^t)\|_1}{T} &\leq \frac{\sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|_1}{T} + \frac{\sum_{t=1}^T \mathbb{E}\|\nabla f_S(\mathbf{w}^t)\|_1}{T} \\ &\leq \frac{\sum_{t=1}^T \sqrt{d}\mathbb{E}\|\nabla f(\mathbf{w}^t) - \nabla f_S(\mathbf{w}^t)\|}{T} + \frac{\sum_{t=1}^T \mathbb{E}\|\nabla f_S(\mathbf{w}^t)\|_1}{T} \\ &= \mathcal{O}\left(\frac{\sqrt{d}}{(1-\delta)T^{1/4}} + \sqrt{\frac{T}{n}}\right). \end{aligned}$$

Based on these two cases, we then prove the proof.

## F. Proofs of Technical Lemmas

### F.1. Proof of Lemma 2

**I.** ( $1 - \alpha \leq \min\left\{\gamma^2, \frac{1}{2T}\right\}$ ) The general case is almost identical to the SGD results proved by [25]. We present the proof here for completeness. We can assume, without loss of generality, that  $S = \{\xi_1, \dots, \xi_{n-1}, \xi_n\}$  and  $S' = \{\xi_1, \dots, \xi_{n-1}, \xi'_n\}$ . Let  $I := \{i_1, \dots, i_T\}$  denote the set of indices selected during the implementation of SignSGDW, and  $(\mathbf{x}^t)_{t \geq 1}$  and  $(\mathbf{y}^t)_{t \geq 1}$  be the sequences generated by the two neighbor sets. By direct computation, we have

$$\begin{aligned} &\mathbb{E}_{\text{SGDW}}\|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \\ &= \mathbb{E}_{\text{SGDW}}\left[\|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \in I\right] \text{Prob}(n \in I) \\ &\quad + \mathbb{E}_{\text{SGDW}}\left[\|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \notin I\right] \text{Prob}(n \notin I). \end{aligned}$$

Noticing that  $\mathbf{y}^T = \mathbf{x}^T$  as  $n \notin I$ , we then derive the following bound

$$\begin{aligned} &\mathbb{E}_{\text{SignSGDW}}\|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \\ &= \mathbb{E}_{\text{SignSGDW}}\left[\|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \in I\right] \text{Prob}(n \in I) \\ &\leq 4G^2[\sum_{t=1}^T \text{Prob}(i_t = n)] \leq \frac{4G^2T}{n} = \mathcal{O}\left(\frac{T}{n}\right). \end{aligned}$$

**II.** ( $1 - \alpha = \gamma$ ) Without loss of generality, we assume that  $S = \{\xi_1, \dots, \xi_{n-1}, \xi_n\}$  and  $S' = \{\xi_1, \dots, \xi_{n-1}, \xi'_n\}$ . Let  $(\mathbf{x}^t)_{t \geq 0}$  and  $(\mathbf{y}^t)_{t \geq 0}$  be the output of SGDW on these two neighbor sets.

1. With probability  $1 - \frac{1}{n}$ , we derive

$$\begin{aligned} \|\mathbf{x}^{t+1} - \mathbf{y}^{t+1}\|^2 &= \|\alpha\mathbf{x}^t - \alpha\mathbf{y}^t - \gamma(\nabla f(\mathbf{x}^t; \xi_{i_t}) - \nabla f(\mathbf{y}^t; \xi_{i_t}))\|^2 \\ &\leq (1+h)\alpha^2\|\mathbf{x}^t - \mathbf{y}^t\|^2 + (1+\frac{1}{h})\gamma^2\|\nabla f(\mathbf{x}^t; \xi_{i_t}) - \nabla f(\mathbf{y}^t; \xi'_{i_t})\|^2. \end{aligned}$$

Setting  $h = \frac{1-\alpha}{\alpha}$ , we are led to

$$\|\mathbf{x}^{t+1} - \mathbf{y}^{t+1}\|^2 \leq \alpha\|\mathbf{x}^t - \mathbf{y}^t\|^2 + \frac{\gamma^2 L^2}{1-\alpha}\|\mathbf{x}^t - \mathbf{y}^t\|^2 \leq \|\mathbf{x}^t - \mathbf{y}^t\|^2.$$

2. With probability  $\frac{1}{n}$ , we get

$$\begin{aligned} \|\mathbf{x}^{t+1} - \mathbf{y}^{t+1}\|^2 &= \|\alpha\mathbf{x}^t - \alpha\mathbf{y}^t - \gamma(\nabla f(\mathbf{x}^t; \xi_n) - \nabla f(\mathbf{y}^t; \xi'_n))\|^2 \\ &\leq \alpha\|\mathbf{x}^t - \mathbf{y}^t\|^2 + \frac{4\gamma^2 G^2}{1-\alpha} \\ &\leq \alpha\|\mathbf{x}^t - \mathbf{y}^t\|^2 + 4\gamma G^2. \end{aligned}$$

In summary, we then get

$$\mathbb{E}\|\mathbf{x}^{t+1} - \mathbf{y}^{t+1}\|^2 \leq \mathbb{E}\|\mathbf{x}^t - \mathbf{y}^t\|^2 + \frac{4\gamma G^2}{n},$$

which then yields

$$\mathbb{E}\|\mathbf{x}^T - \mathbf{y}^T\|^2 \leq \frac{4\gamma T G^2}{n}.$$

The Lipschitz property of the gradient then indicates

$$\sup_{\xi \in \Xi} \mathbb{E}\|\nabla f(\mathbf{x}^T; \xi) - \nabla f(\mathbf{y}^T; \xi)\|^2 \leq L^2 \mathbb{E}\|\mathbf{x}^T - \mathbf{y}^T\|^2 \leq \frac{4\gamma T L^2 G^2}{n}.$$

Hence, we can get the bound as

$$\sup_{\xi \in \Xi} \mathbb{E}\|\nabla f(\mathbf{x}^T; \xi) - \nabla f(\mathbf{y}^T; \xi)\|^2 \leq \frac{4\gamma T L^2 G^2}{n}.$$

## F.2. Proof of Lemma 3

Similarly, we can get

$$\begin{aligned} \mathbb{E}\|\mathbf{w}^{t+1}\|^2 &= \mathbb{E}\|\alpha\mathbf{w}^t - \gamma\mathbf{g}^t\|^2 \\ &= \mathbb{E}\|\alpha\mathbf{w}^t\|^2 - 2\gamma\alpha\langle\mathbf{w}^t, \mathbf{g}^t\rangle + \gamma^2\mathbb{E}\|\mathbf{g}^t\|^2 \\ &= \mathbb{E}\|\alpha\mathbf{w}^t\|^2 - 2\gamma\alpha\mathbb{E}\langle\mathbf{w}^t, \nabla f(\mathbf{w}^t)\rangle + \gamma^2\mathbb{E}\|\nabla f(\mathbf{w}^t)\|^2 + \gamma^2\sigma^2 \\ &= \mathbb{E}\|\alpha\mathbf{w}^t\|^2 + 2\gamma\alpha\mathbb{E}(\|\mathbf{w}^t\| \cdot \|\nabla f(\mathbf{w}^t)\|) + \gamma^2\mathbb{E}\|\nabla f(\mathbf{w}^t)\|^2 + \gamma^2\sigma^2 \\ &\leq (1 + \frac{1}{h})\alpha^2\mathbb{E}\|\mathbf{w}^t\|^2 + (1 + h)\gamma^2\mathbb{E}\|\nabla f(\mathbf{w}^t)\|^2 + \gamma^2\sigma^2. \end{aligned}$$

Setting  $h = \frac{\alpha}{1-\alpha}$ , we get

$$\mathbb{E}\|\mathbf{w}^{t+1}\|^2 \leq \alpha\mathbb{E}\|\mathbf{w}^t\|^2 + \frac{\gamma^2}{1-\alpha}\mathbb{E}\|\nabla f(\mathbf{w}^t)\|^2 + \gamma^2\sigma^2.$$

Thus, we are led to

$$\begin{aligned} \mathbb{E}\|\mathbf{w}^{t+1}\|^2 &\leq \alpha^t\|\mathbf{w}^1\|^2 + \sum_{j=1}^t \alpha^{t-j}(\frac{1}{1-\alpha}\gamma^2\mathbb{E}\|\nabla f(\mathbf{w}^j)\|^2 + \gamma^2\sigma^2) \\ &\leq \alpha^t\|\mathbf{w}^1\|^2 + \frac{1}{1-\alpha}\gamma^2 \sum_{j=1}^t \alpha^{t-j}\mathbb{E}\|\nabla f(\mathbf{w}^j)\|^2 + \frac{\gamma^2(1-\alpha^T)}{1-\alpha}\sigma^2 \end{aligned}$$

as  $t \leq T$ . Summation from  $t = 1$  to  $T$  leads to

$$\sum_{t=1}^T \mathbb{E}\|\mathbf{w}^t\|^2 \leq \sum_{t=1}^T \mathbb{E}\|\mathbf{w}^{t+1}\|^2 \leq \frac{1}{1-\alpha}\|\mathbf{w}^1\|^2 + \frac{\gamma^2(1-\alpha^T)}{(1-\alpha)^2} \sum_{j=1}^T \mathbb{E}\|\nabla f(\mathbf{w}^j)\|^2 + \frac{\gamma^2 T (1-\alpha^T)}{1-\alpha} \sigma^2.$$

Noting that  $\mathbf{w}^1 = \mathbf{0}$ , we then proved the result.

## F.3. Proof of Lemma 4

With the scheme of the algorithm, we have

$$\|\mathbf{w}^{t+1}\| \leq \alpha\|\mathbf{w}^t\| + \gamma\sqrt{d} \leq \gamma\sqrt{d}(\sum_{j=0}^t \alpha^j) \leq \gamma\sqrt{d}(\sum_{j=0}^{T-1} \alpha^j) \leq \frac{1-\alpha^T}{1-\alpha}\gamma\sqrt{d}.$$

For any  $i \in \{1, 2, \dots, n\}$ , it follows that

$$\|\mathbf{w}_i^{t+1}\| \leq \alpha\|\mathbf{w}_i^t\| + \gamma \leq \gamma(\sum_{j=0}^t \alpha^j) \leq \gamma(\sum_{j=0}^{T-1} \alpha^j) \leq \frac{1-\alpha^T}{1-\alpha}\gamma.$$

#### F.4. Proof of Lemma 5

The proof of this lemma is inspired by the work of [34]. With the smoothness of  $\nabla f_S$ , we have

$$\begin{aligned}
& f_S(\mathbf{x}^\dagger) - f_S(\mathbf{x}^\dagger) \\
& \leq \langle \nabla f_S(\mathbf{x}^\dagger), \mathbf{x}^\dagger - \mathbf{x}^\dagger \rangle + \frac{L}{2} \|\mathbf{x}^\dagger - \mathbf{x}^\dagger\|^2 \\
& \leq -\gamma \langle \nabla f(\mathbf{x}^\dagger), \text{Sign}(\mathbf{m}) \rangle + (\alpha - 1) \langle \nabla f_S(\mathbf{x}^\dagger), \mathbf{x}^\dagger \rangle + L \|\gamma \text{Sign}(\mathbf{m})\|^2 + L \|(\alpha - 1)\mathbf{x}^\dagger\|^2 \\
& \leq -\gamma \langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}(\mathbf{m}) \rangle + Ld\gamma^2 + L(\alpha - 1)^2 \|\mathbf{x}^\dagger\|^2 + (1 - \alpha)\delta \|\nabla f_S(\mathbf{x}^\dagger)\|_1 \cdot \|\mathbf{x}^\dagger\|_\infty \\
& = -\gamma \langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}[\nabla f_S(\mathbf{x}^\dagger)] \rangle + \gamma \langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}(\mathbf{m}) - \text{Sign}[\nabla f_S(\mathbf{x}^\dagger)] \rangle \\
& \quad + Ld\gamma^2 + L(\alpha - 1)^2 \|\mathbf{x}^\dagger\|^2 + (1 - \alpha)\delta \|\nabla f_S(\mathbf{x}^\dagger)\|_1 \cdot \|\mathbf{x}^\dagger\|_\infty \\
& = -\gamma \|\nabla f_S(\mathbf{x}^\dagger)\|_1 + \gamma \langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}(\mathbf{m}) - \text{Sign}[\nabla f_S(\mathbf{x}^\dagger)] \rangle \\
& \quad + Ld\gamma^2 + L(\alpha - 1)^2 \|\mathbf{x}^\dagger\|^2 + (1 - \alpha)\delta \|\nabla f_S(\mathbf{x}^\dagger)\|_1 \cdot \|\mathbf{x}^\dagger\|_\infty.
\end{aligned}$$

With direct computation,

$$\langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}[\nabla f_S(\mathbf{x}^\dagger)] - \text{Sign}(\mathbf{m}) \rangle = \sum_{i=1}^d [\nabla f_S(\mathbf{x}^\dagger)]_i \cdot [\text{Sign}([\nabla f_S(\mathbf{x}^\dagger)]_i) - \text{Sign}([\mathbf{m}]_i)].$$

**I.** If  $\text{Sign}([\nabla f_S(\mathbf{x}^\dagger)]_i) = \text{Sign}([\mathbf{m}]_i)$ ,  $[\nabla f_S(\mathbf{x}^\dagger)]_i \cdot [\text{Sign}([\nabla f_S(\mathbf{x}^\dagger)]_i) - \text{Sign}([\mathbf{m}]_i)] = 0$ .

**II.** Otherwise,  $[\nabla f_S(\mathbf{x}^\dagger)]_i \cdot [\mathbf{m}]_i \leq 0$  gives us

$$|\epsilon_i| = |[\nabla f_S(\mathbf{x}^\dagger)]_i - [\mathbf{m}]_i| \geq |[\nabla f_S(\mathbf{x}^\dagger)]_i|.$$

And we can get

$$[\nabla f_S(\mathbf{x}^\dagger)]_i \cdot [\text{Sign}([\nabla f_S(\mathbf{x}^\dagger)]_i) - \text{Sign}([\mathbf{m}]_i)] \leq 2|[\nabla f_S(\mathbf{x}^\dagger)]_i| \leq 2|\epsilon_i|. \quad (9)$$

Therefore, for any  $i \in [d]$ , inequality (9) holds. Based on this, we have

$$\langle \nabla f_S(\mathbf{x}^\dagger), \text{Sign}[\nabla f_S(\mathbf{x}^\dagger)] - \text{Sign}(\mathbf{m}) \rangle \leq 2 \sum_{i=1}^d |\epsilon_i| \leq 2\sqrt{d}\|\epsilon\|,$$

which further yields

$$\begin{aligned}
f_S(\mathbf{x}^\dagger) - f_S(\mathbf{x}^\dagger) & \leq -\gamma \|\nabla f(\mathbf{x}^\dagger)\|_1 + 2\sqrt{d}\gamma\|\epsilon\| + Ld\gamma^2 \\
& \quad + L(\alpha - 1)^2 \|\mathbf{x}^\dagger\|^2 + (1 - \alpha)\delta \|\nabla f_S(\mathbf{x}^\dagger)\|_1 \cdot \|\mathbf{x}^\dagger\|_\infty.
\end{aligned}$$

Recall that  $\|\mathbf{x}^\dagger\| \leq \frac{\gamma\sqrt{d}}{1-\alpha}$ ,  $\|\mathbf{x}^\dagger\|_\infty \leq \frac{1-\alpha^T}{1-\alpha}\gamma$ , we then get

$$\begin{aligned}
f_S(\mathbf{x}^\dagger) - f_S(\mathbf{x}^\dagger) & \leq -\gamma \|\nabla f_S(\mathbf{x}^\dagger)\|_1 + 2\sqrt{d}\gamma\|\epsilon\| + 2Ld\gamma^2 + (1 - \alpha^T)\delta\gamma \|\nabla f(\mathbf{x}^\dagger)\|_1 \\
& = -\gamma[1 - (1 - \alpha^T)\delta] \|\nabla f(\mathbf{x}^\dagger)\|_1 + 2\sqrt{d}\gamma\|\epsilon\| + 2Ld\gamma^2.
\end{aligned}$$

#### F.5. Proof of Lemma 6

Without loss of generality, we assume that  $S = \{\xi_1, \dots, \xi_{n-1}, \xi_n\}$  and  $S' = \{\xi_1, \dots, \xi_{n-1}, \xi'_n\}$ . Let  $I := \{i_1, \dots, i_T\}$  be the set of indices selected during the implementation of SignSGDW, and  $(\mathbf{x}^t)_{t \geq 1}$  and  $(\mathbf{y}^t)_{t \geq 1}$  be generated by the two neighbor sets. The scheme of the algorithm gives us

$$\begin{aligned}
\|\mathbf{x}^{t+1} - \mathbf{y}^{t+1}\| & = \|\alpha\mathbf{x}^t - \alpha\mathbf{y}^t - \gamma(\text{Sign}(\mathbf{m}^t) - \text{Sign}(\hat{\mathbf{m}}^t))\| \\
& \leq \alpha \|\mathbf{x}^t - \mathbf{y}^t\|^2 + \gamma \|\text{Sign}(\mathbf{m}^t) - \text{Sign}(\hat{\mathbf{m}}^t)\| \\
& \leq \alpha \|\mathbf{x}^t - \mathbf{y}^t\| + 2\gamma\sqrt{d}.
\end{aligned}$$

Hence, we are led to

$$\sup_{t \geq 1} \|\mathbf{x}^t - \mathbf{y}^t\| \leq \frac{2\gamma\sqrt{d}(1-\alpha^T)}{(1-\alpha)}.$$

It then follows

$$\begin{aligned} & \mathbb{E}_{\text{SignSGDW}} \|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \\ &= \mathbb{E}_{\text{SignSGDW}} \left[ \|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \in I \right] \text{Prob}(n \in I) \\ &\quad + \mathbb{E}_{\text{SignSGDW}} \left[ \|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \notin I \right] \text{Prob}(n \notin I). \end{aligned}$$

Noticing that  $\mathbf{y}^T = \mathbf{x}^T$  as  $n \notin I$ , we then get

$$\begin{aligned} & \mathbb{E}_{\text{SignSGDW}} \|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \\ &= \mathbb{E}_{\text{SignSGDW}} \left[ \|\nabla f(\mathbf{y}^T; \xi) - \nabla f(\mathbf{x}^T; \xi)\|^2 \middle| n \in I \right] \text{Prob}(n \in I) \\ &\leq L^2 \mathbb{E} \|\mathbf{x}^T - \mathbf{y}^T\|^2 \text{Prob}(n \in I) \leq \frac{4L^2\gamma^2 d(1-\alpha^T)^2}{(1-\alpha)^2} \text{Prob}(n \in I) \\ &\leq \frac{4L^2\gamma^2 d(1-\alpha^T)^2}{(1-\alpha)^2} [\sum_{t=1}^T \text{Prob}(i_t = n)] \leq \frac{4L^2\gamma^2 d(1-\alpha^T)^2}{(1-\alpha)^2} \cdot \frac{T}{n}. \end{aligned} \tag{10}$$

**I.** In the general case ( $\delta = 1$ ),  $1 - \alpha \leq \frac{1}{T}$ , it holds  $\frac{1-\alpha^T}{1-\alpha} \asymp T$  yielding

$$\frac{4L^2\gamma^2 d(1-\alpha^T)^2}{(1-\alpha)^2} \cdot \frac{T}{n} = \mathcal{O}\left(\frac{4L^2\gamma^2 dT^3}{n}\right) = \mathcal{O}\left(\frac{T^{3/2}}{n}\right).$$

**II.** If  $\delta < 1$  and  $\alpha = 1 - \frac{1}{T^{1/4}}$ , we have

$$\frac{4L^2\gamma^2 d(1-\alpha^T)^2}{(1-\alpha)^2} \cdot \frac{T}{n} \leq \frac{4L^2\gamma^2 d}{(1-\alpha)^2} \cdot \frac{T}{n} = \mathcal{O}\left(\frac{4L^2\gamma^2 dT^{5/2}}{n}\right) = \mathcal{O}\left(\frac{T}{n}\right).$$