

A Extra Lemmas

In this section we (re)state and prove some Lemmas.

First, we provide the proof of Lemma 1, restated below for convenience.

Lemma 1. *Suppose $\eta_t \leq \frac{1}{4L}$ for all t . Then*

$$\mathbb{E}[F(\mathbf{x}_{t+1}) - F(\mathbf{x}_t)] \leq \mathbb{E} \left[-\eta_t/4 \|\nabla F(\mathbf{x}_t)\|^2 + 3\eta_t/4 \|\boldsymbol{\epsilon}_t\|^2 \right].$$

Proof. Using the smoothness of F and the definition of \mathbf{x}_{t+1} from the algorithm, we have

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}_{t+1})] &\leq \mathbb{E} \left[F(\mathbf{x}_t) - \nabla F(\mathbf{x}_t) \cdot \eta_t \mathbf{d}_t + \frac{L\eta_t^2}{2} \|\mathbf{d}_t\|^2 \right] \\ &= \mathbb{E} \left[F(\mathbf{x}_t) - \eta_t \|\nabla F(\mathbf{x}_t)\|^2 - \eta_t \nabla F(\mathbf{x}_t) \cdot \boldsymbol{\epsilon}_t + \frac{L\eta_t^2}{2} \|\mathbf{d}_t\|^2 \right] \\ &\leq \mathbb{E} \left[F(\mathbf{x}_t) - \frac{\eta_t}{2} \|\nabla F(\mathbf{x}_t)\|^2 + \frac{\eta_t}{2} \|\boldsymbol{\epsilon}_t\|^2 + \frac{L\eta_t^2}{2} \|\mathbf{d}_t\|^2 \right] \\ &\leq \mathbb{E} \left[F(\mathbf{x}_t) - \frac{\eta_t}{2} \|\nabla F(\mathbf{x}_t)\|^2 + \frac{\eta_t}{2} \|\boldsymbol{\epsilon}_t\|^2 + L\eta_t^2 \|\boldsymbol{\epsilon}_t\|^2 + L\eta_t^2 \|\nabla F(\mathbf{x}_t)\|^2 \right] \\ &\leq \mathbb{E} \left[F(\mathbf{x}_t) - \frac{\eta_t}{2} \|\nabla F(\mathbf{x}_t)\|^2 + \frac{3\eta_t}{4} \|\boldsymbol{\epsilon}_t\|^2 + \frac{\eta_t}{4} \|\nabla F(\mathbf{x}_t)\|^2 \right], \end{aligned}$$

where in the second inequality we used Young's inequality, the third one uses $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2$, and the last one uses $\eta_t \leq 1/4L$. \square

This next Lemma is a technical observation that is important for the proof of Lemma 2.

Lemma 3.

$$\begin{aligned} \mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \right] &= 0 \\ \mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - \nabla F(\mathbf{x}_t) + \nabla F(\mathbf{x}_{t-1})) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \right] &= 0. \end{aligned}$$

Proof. From inspection of the update formula, the hypothesis implies that $\boldsymbol{\epsilon}_{t-1} = \mathbf{d}_{t-1} - \nabla F(\mathbf{x}_{t-1})$ and \mathbf{x}_t are both independent of ξ_t . Then, by first taking expectation with respect to ξ_t and then with respect to ξ_1, \dots, ξ_{t-1} , we obtain

$$\begin{aligned} \mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \right] \\ = \mathbb{E} \left[\mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \mid \xi_1, \dots, \xi_{t-1} \right] \right] = 0. \end{aligned}$$

Analogously, for the second equality we have

$$\begin{aligned} \mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - \nabla F(\mathbf{x}_t) + \nabla F(\mathbf{x}_{t-1})) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \right] \\ = \mathbb{E} \left[\mathbb{E} \left[(\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - (\nabla F(\mathbf{x}_t) - \nabla F(\mathbf{x}_{t-1}))) \cdot \eta_{t-1}^{-1} (1 - a_t)^2 \boldsymbol{\epsilon}_{t-1} \mid \xi_1, \dots, \xi_{t-1} \right] \right] \\ = 0. \quad \square \end{aligned}$$

The following Lemma is a standard consequence of convexity.

Lemma 4. *Let $a_0 > 0$ and $a_1, \dots, a_T \geq 0$. Then*

$$\sum_{t=1}^T \frac{a_t}{a_0 + \sum_{i=1}^t a_i} \leq \ln \left(1 + \frac{\sum_{i=1}^T a_i}{a_0} \right).$$

Proof. By the concavity of the log function, we have

$$\ln \left(a_0 + \sum_{i=1}^t a_i \right) - \ln \left(a_0 + \sum_{i=1}^{t-1} a_i \right) \geq \frac{a_t}{a_0 + \sum_{i=1}^t a_i}.$$

Summing over $t = 1, \dots, T$ both sides of the inequality, we have the stated bound. \square

A.1 Proof of Lemma 2

In this section we present the deferred proof of Lemma 2, restating the result below for reference

Lemma 2. *With the notation in Algorithm 1, we have*

$$\mathbb{E} [\|\epsilon_t\|^2/\eta_{t-1}] \leq \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + (1-a_t)^2(1+4L^2\eta_{t-1}^2)\|\epsilon_{t-1}\|^2/\eta_{t-1} + 4(1-a_t)^2L^2\eta_{t-1}\|\nabla F(\mathbf{x}_{t-1})\|^2] .$$

Proof. First, observe that

$$\begin{aligned} & \mathbb{E} [\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)\|^2] \\ &= \mathbb{E} [\eta_{t-1}^3(\|\nabla f(\mathbf{x}_t, \xi_t)\|^2 + \|\nabla F(\mathbf{x}_t)\|^2 - 2\nabla f(\mathbf{x}_t, \xi_t) \cdot \nabla F(\mathbf{x}_t))] \\ &= \mathbb{E} [\eta_{t-1}^3\mathbb{E} [\|\nabla f(\mathbf{x}_t, \xi_t)\|^2 + \|\nabla F(\mathbf{x}_t)\|^2 - 2\nabla f(\mathbf{x}_t, \xi_t) \cdot \nabla F(\mathbf{x}_t) | \xi_1, \dots, \xi_{t-1}]] \\ &= \mathbb{E} [\eta_{t-1}^3(\|\nabla f(\mathbf{x}_t, \xi_t)\|^2 - \|\nabla F(\mathbf{x}_t)\|^2)] \\ &\leq \mathbb{E} [\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t)\|^2] . \end{aligned} \tag{5}$$

In the same way, we also have that

$$\begin{aligned} & \mathbb{E} [\eta_{t-1}^{-1}(1-a_t^2)\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - \nabla F(\mathbf{x}_t) + \nabla F(\mathbf{x}_{t-1})\|^2] \\ &\leq \mathbb{E} [\eta_{t-1}^{-1}(1-a_t^2)\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t)\|^2] . \end{aligned} \tag{6}$$

By definition of ϵ_t and the notation in Algorithm 1, we have $\epsilon_t = \mathbf{d}_t - \nabla F(\mathbf{x}_t) = \nabla f(\mathbf{x}_t, \xi_t) + (1-a_t)(\mathbf{d}_{t-1} - \nabla f(\mathbf{x}_{t-1}, \xi_t)) - \nabla F(\mathbf{x}_t)$. Hence, we can write

$$\begin{aligned} \mathbb{E} [\eta_{t-1}^{-1}\|\epsilon_t\|^2] &= \mathbb{E} [\eta_{t-1}^{-1}\|\nabla f(\mathbf{x}_t, \xi_t) + (1-a_t)(\mathbf{d}_{t-1} - \nabla f(\mathbf{x}_{t-1}, \xi_t)) - \nabla F(\mathbf{x}_t)\|^2] \\ &= \mathbb{E} [\eta_{t-1}^{-1}\|a_t(\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)) + (1-a_t)(\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - \nabla F(\mathbf{x}_t) + \nabla F(\mathbf{x}_{t-1})) \\ &\quad + (1-a_t)(\mathbf{d}_{t-1} - \nabla F(\mathbf{x}_{t-1}))\|^2] \\ &\leq \mathbb{E} [2c^2\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)\|^2 + 2\eta_{t-1}^{-1}(1-a_t)^2\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t) - \nabla F(\mathbf{x}_t) + \nabla F(\mathbf{x}_{t-1})\|^2 \\ &\quad + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &\leq \mathbb{E} [2c^2\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t)\|^2 + 2\eta_{t-1}^{-1}(1-a_t)^2\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla f(\mathbf{x}_{t-1}, \xi_t)\|^2 + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &\leq \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + 2\eta_{t-1}^{-1}(1-a_t)^2L^2\|\mathbf{x}_t - \mathbf{x}_{t-1}\|^2 + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &= \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + 2(1-a_t)^2L^2\eta_{t-1}\|\mathbf{d}_{t-1}\|^2 + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &= \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + 2(1-a_t)^2L^2\eta_{t-1}\|\epsilon_{t-1} + \nabla F(\mathbf{x}_{t-1})\|^2 + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &\leq \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + 4(1-a_t)^2L^2\eta_{t-1}(\|\epsilon_{t-1}\|^2 + \|\nabla F(\mathbf{x}_{t-1})\|^2) + \eta_{t-1}^{-1}(1-a_t)^2\|\epsilon_{t-1}\|^2] \\ &= \mathbb{E} [2c^2\eta_{t-1}^3G_t^2 + \eta_{t-1}^{-1}(1-a_t)^2(1+4L^2\eta_{t-1}^2)\|\epsilon_{t-1}\|^2 + 4(1-a_t)^2L^2\eta_{t-1}\|\nabla F(\mathbf{x}_{t-1})\|^2] , \end{aligned}$$

where in the first inequality we used Lemma 3 (See Appendix A) and $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2$, in the second inequality we used (5) and (6), in the third one the Lipschitzness and smoothness of the functions f , and in the last inequality we used again $\|\mathbf{x} + \mathbf{y}\|^2 \leq 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2$. \square

B Non-adaptive Bound Without Lipschitz Assumption

In our analysis of STORM in Theorem 1 we assume that the losses are G -Lipschitz for some known constant G with probability 1. Often this kind of Lipschitz assumption is avoided in other variance-reduction analyses [18, 8, 25]. These works also require oracle knowledge of the parameter σ . It turns out that our use of this assumption is actually only necessary in order to facilitate our adaptive analysis - in fact even for ordinary (non-variance-reduced) gradient descent methods the Lipschitz assumption seems to be a common thread in adaptive analyses [16, 28]. If we are given access to the true value of σ , then we can choose a deterministic learning rate schedule in order to avoid requiring a Lipschitz bound. All that needs be done is replace all instances of G or G_t in STORM with the oracle-tuned value σ , which we outline in Algorithm 2 below.

The convergence guarantee of Algorithm 2 is presented in Theorem 2 below, which is nearly identical to Theorem 1 but losses adaptivity to σ in exchange for removing the G -Lipschitz requirement.

Algorithm 2 STORM without Lipschitz Bound

- 1: **Input:** Parameters k, w, c , initial point \mathbf{x}_1
 - 2: **Sample** ξ_1
 - 3: $G_1 \leftarrow \|\nabla f(\mathbf{x}_1, \xi_1)\|$
 - 4: $\mathbf{d}_1 \leftarrow \nabla f(\mathbf{x}_1, \xi_1)$
 - 5: $\eta_0 \leftarrow \frac{k}{w^{1/3}}$
 - 6: **for** $t = 1$ **to** T **do**
 - 7: $\eta_t \leftarrow \frac{k}{(w + \sigma^2 t)^{1/3}}$
 - 8: $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta_t \mathbf{d}_t$
 - 9: $a_{t+1} \leftarrow c\eta_t^2$
 - 10: **Sample** ξ_{t+1}
 - 11: $G_{t+1} \leftarrow \|\nabla f(\mathbf{x}_{t+1}, \xi_{t+1})\|$
 - 12: $\mathbf{d}_{t+1} \leftarrow \nabla f(\mathbf{x}_{t+1}, \xi_{t+1}) + (1 - a_{t+1})(\mathbf{d}_t - \nabla f(\mathbf{x}_t, \xi_{t+1}))$
 - 13: **end for**
 - 14: Choose $\hat{\mathbf{x}}$ uniformly at random from $\mathbf{x}_1, \dots, \mathbf{x}_T$. (In practice, set $\hat{\mathbf{x}} = \mathbf{x}_T$).
 - 15: **return** $\hat{\mathbf{x}}$
-

Theorem 2. Under the assumptions in Section 3, for any $b > 0$, we write $k = \frac{b\sigma^{3/2}}{L}$. Set $c = 28L^2 + \sigma^2/(7Lk^3) = L^2(28 + 1/(7b^3))$ and $w = \max\left((4Lk)^3, 2\sigma^2, \left(\frac{ck}{4L}\right)^3\right) = \sigma^2 \max\left((4b)^3, 2, (28b + \frac{1}{7b^2})^3/64\right)$. Then, Algorithm 2 satisfies

$$\frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \|\nabla F(\mathbf{x}_t)\|^2 \right] \leq \frac{M \frac{w^{1/3}}{k}}{T} + \frac{M \frac{w\sigma^{2/3}}{k}}{T^{2/3}},$$

where $M = 8(F(\mathbf{x}_1) - F^*) + \frac{w^{1/3}\sigma^2}{4L^2k} + \frac{k^3c^2}{2L^2} \ln(T+2)$.

In order to prove this Theorem, we need a non-adaptive analog of Lemma 2:

Lemma 5. With the notation in Algorithm 2, we have

$$\mathbb{E} \left[\frac{\|\epsilon_t\|^2}{\eta_{t-1}} \right] \leq \mathbb{E} \left[2c^2\eta_{t-1}^3\sigma^2 + \frac{(1-a_t)^2(1+4L^2\eta_{t-1}^2)\|\epsilon_{t-1}\|^2}{\eta_{t-1}} + 4(1-a_t)^2L^2\eta_{t-1}\|\nabla F(\mathbf{x}_{t-1})\|^2 \right].$$

Proof. The proof is nearly identical to that of Lemma 2: the only difference is that instead of using the identity $\mathbb{E}[\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)\|^2] \leq \mathbb{E}[\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t)\|^2] = \mathbb{E}[\eta_{t-1}^3G_t^2]$, we directly use the value of σ : $\mathbb{E}[\eta_{t-1}^3\|\nabla f(\mathbf{x}_t, \xi_t) - \nabla F(\mathbf{x}_t)\|^2] \leq \eta_{t-1}^3\sigma^2$. \square

Now we can prove Theorem 2:

Proof of Theorem 2. This proof is also nearly identical to the analogous adaptive result of Theorem 1.

Again, we consider the potential $\Phi_t = F(\mathbf{x}_t) + \frac{1}{32L^2\eta_{t-1}}\|\epsilon_t\|^2$ and upper bound $\Phi_{t+1} - \Phi_t$ for each t .

Since $w \geq (4Lk)^3$, we have $\eta_t \leq \frac{1}{4L}$. Further, since $a_{t+1} = c\eta_t^2$, we have $a_{t+1} \leq \frac{ck}{4Lw^{1/3}} \leq 1$ for all t . Then, we first consider $\eta_t^{-1}\|\epsilon_{t+1}\|^2 - \eta_{t-1}^{-1}\|\epsilon_t\|^2$. Using Lemma 5, we obtain

$$\begin{aligned} & \mathbb{E} \left[\eta_t^{-1}\|\epsilon_{t+1}\|^2 - \eta_{t-1}^{-1}\|\epsilon_t\|^2 \right] \\ & \leq \mathbb{E} \left[2c^2\eta_t^3\sigma^2 + \frac{(1-a_{t+1})^2(1+4L^2\eta_t^2)\|\epsilon_t\|^2}{\eta_t} + 4(1-a_{t+1})^2L^2\eta_t\|\nabla F(\mathbf{x}_t)\|^2 - \frac{\|\epsilon_t\|^2}{\eta_{t-1}} \right] \\ & \leq \mathbb{E} \left[\underbrace{2c^2\eta_t^3\sigma^2}_{A_t} + \underbrace{\left(\eta_t^{-1}(1-a_{t+1})(1+4L^2\eta_t^2) - \eta_{t-1}^{-1} \right) \|\epsilon_t\|^2}_{B_t} + \underbrace{4L^2\eta_t\|\nabla F(\mathbf{x}_t)\|^2}_{C_t} \right]. \end{aligned}$$

Let us focus on the terms of this expression individually. For the first term, A_t , observe that $w \geq 2\sigma^2$ to obtain:

$$\begin{aligned} \sum_{t=1}^T A_t &= \sum_{t=1}^T 2c^2 \eta_t^3 \sigma^2 = \sum_{t=1}^T \frac{2k^3 c^2 \sigma^2}{w + t\sigma^2} \leq \sum_{t=1}^T \frac{2k^3 c^2}{t+1} \\ &\leq 2k^3 c^2 \ln(T+2). \end{aligned}$$

For the second term B_t , we have

$$B_t \leq (\eta_t^{-1} - \eta_{t-1}^{-1} + \eta_t^{-1}(4L^2 \eta_t^2 - a_{t+1})) \|\epsilon_t\|^2 = (\eta_t^{-1} - \eta_{t-1}^{-1} + \eta_t(4L^2 - c)) \|\epsilon_t\|^2.$$

Let us focus on $\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}}$ for a minute. Using the concavity of $x^{1/3}$, we have $(x+y)^{1/3} \leq x^{1/3} + yx^{-2/3}/3$. Therefore:

$$\begin{aligned} \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} &= \frac{1}{k} \left[(w + t\sigma^2)^{1/3} - (w + (t-1)\sigma^2)^{1/3} \right] \leq \frac{\sigma^2}{3k(w + (t-1)\sigma^2)^{2/3}} \\ &\leq \frac{\sigma^2}{3k(w - \sigma^2 + t\sigma^2)^{2/3}} \leq \frac{\sigma^2}{3k(w/2 + t\sigma^2)^{2/3}} \\ &\leq \frac{2^{2/3}\sigma^2}{3k(w + t\sigma^2)^{2/3}} \leq \frac{2^{2/3}\sigma^2}{3k^3} \eta_t^2 \leq \frac{2^{2/3}\sigma^2}{12Lk^3} \eta_t \leq \frac{\sigma^2}{7Lk^3} \eta_t, \end{aligned}$$

where we have used that that $w \geq (4Lk)^3$ to have $\eta_t \leq \frac{1}{4L}$.

Further, since $c = 28L^2 + \sigma^2/(7Lk^3)$, we have

$$\eta_t(4L^2 - c) \leq -24L^2 \eta_t - \sigma^2 \eta_t / (7Lk^3).$$

Thus, we obtain $B_t \leq -24L^2 \eta_t \|\epsilon_t\|^2$. Putting all this together yields:

$$\frac{1}{32L^2} \sum_{t=1}^T \left(\frac{\|\epsilon_{t+1}\|^2}{\eta_t} - \frac{\|\epsilon_t\|^2}{\eta_{t-1}} \right) \leq \frac{k^3 c^2}{16L^2} \ln(T+2) + \sum_{t=1}^T \left[\frac{\eta_t}{8} \|\nabla F(\mathbf{x}_t)\|^2 - \frac{3\eta_t}{4} \|\epsilon_t\|^2 \right]. \quad (7)$$

Now, we analyze the potential Φ_t . This analysis is completely identical to that of Theorem 1, and is only reproduced here for convenience. Since $\eta_t \leq \frac{1}{4L}$, we can use Lemma 1 to obtain

$$\mathbb{E}[\Phi_{t+1} - \Phi_t] \leq \mathbb{E} \left[-\frac{\eta_t}{4} \|\nabla F(\mathbf{x}_t)\|^2 + \frac{3\eta_t}{4} \|\epsilon_t\|^2 + \frac{1}{32L^2 \eta_t} \|\epsilon_{t+1}\|^2 - \frac{1}{32L^2 \eta_{t-1}} \|\epsilon_t\|^2 \right].$$

Summing over t and using (7), we obtain

$$\begin{aligned} \mathbb{E}[\Phi_{T+1} - \Phi_1] &\leq \sum_{t=1}^T \mathbb{E} \left[-\frac{\eta_t}{4} \|\nabla F(\mathbf{x}_t)\|^2 + \frac{3\eta_t}{4} \|\epsilon_t\|^2 + \frac{1}{32L^2 \eta_t} \|\epsilon_{t+1}\|^2 - \frac{1}{32L^2 \eta_{t-1}} \|\epsilon_t\|^2 \right] \\ &\leq \mathbb{E} \left[\frac{k^3 c^2}{16L^2} \ln(T+2) - \sum_{t=1}^T \frac{\eta_t}{8} \|\nabla F(\mathbf{x}_t)\|^2 \right]. \end{aligned}$$

Reordering the terms, we have

$$\begin{aligned} \mathbb{E} \left[\sum_{t=1}^T \eta_t \|\nabla F(\mathbf{x}_t)\|^2 \right] &\leq \mathbb{E} \left[8(\Phi_1 - \Phi_{T+1}) + \frac{k^3 c^2}{2L^2} \ln(T+2) \right] \\ &\leq 8(F(\mathbf{x}_1) - F^*) + \frac{1}{4L^2 \eta_0} \mathbb{E}[\|\epsilon_1\|^2] + \frac{k^3 c^2}{2L^2} \ln(T+2) \\ &\leq 8(F(\mathbf{x}_1) - F^*) + \frac{w^{1/3} \sigma^2}{4L^2 k} + \frac{k^3 c^2}{2L^2} \ln(T+2), \end{aligned}$$

where the last inequality is given by the definition of d_1 and η_0 in the algorithm.

At this point the rest of the proof could proceed in an identical manner to that of Theorem 1. However, since η_t is now independent of $\nabla F(x_t)$ by virtue of being deterministic, we can simplify the remainder of the proof somewhat by avoiding the use of Cauchy-Schwarz inequality.

Since η_t is deterministic, we have $\mathbb{E} \left[\sum_{t=1}^T \eta_t \|\nabla F(\mathbf{x}_t)\|^2 \right] \geq \eta_T \mathbb{E} \left[\sum_{t=1}^T \|\nabla F(\mathbf{x}_t)\|^2 \right]$. Then divide by $T\eta_T$ to conclude

$$\frac{1}{T} \mathbb{E} \left[\sum_{t=1}^T \|\nabla F(\mathbf{x}_t)\|^2 \right] \leq \frac{M \frac{w^{1/3}}{k}}{T} + \frac{M \frac{w\sigma^{2/3}}{k}}{T^{2/3}},$$

where we have used the definition $M = 8(F(\mathbf{x}_1) - F^*) + \frac{w^{1/3}\sigma^2}{4L^2k} + \frac{k^3c^2}{2L^2} \ln(T+2)$ and the identity $(a+b)^{1/3} \leq a^{1/3} + b^{1/3}$ □