

SUPPLEMENTAL MATERIALS

Lemma 5. [15] If $f(\mathbf{x})$ is μ -strongly convex and \mathbf{x}_* denotes the optimal solution to $\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$. For any $\mathbf{x} \in \mathcal{D}$, we have $f(\mathbf{x}) - f(\mathbf{x}_*) \leq 2G_1^2/\mu$.

Proof. From Assumption A1, we have $\|\nabla f(\mathbf{x})\|_2 \leq G_1$. Hence

$$f(\mathbf{x}) - f(\mathbf{x}_*) \leq G_1 \|\mathbf{x} - \mathbf{x}_*\|_2.$$

Moreover from the strong convexity in $f(\cdot)$ we have

$$f(\mathbf{x}) - f(\mathbf{x}_*) \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2.$$

From the two inequalities above, we can easily verify that

$$\|\mathbf{x} - \mathbf{x}_*\|_2 \leq \frac{2G_1}{\mu}, \quad f(\mathbf{x}) - f(\mathbf{x}_*) \leq \frac{2G_1^2}{\mu}.$$

This completes the proof. \square

Proof of Theorem 2

The proof of Theorem 2 is based on an important result, as summarized in Lemma 6.

Lemma 6. [20] Assume $\|\mathbf{x}_* - \mathbf{x}_t\|_2 \leq D$ for all t . Define $D_T = \frac{1}{T} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|_2^2$ and $\Lambda_T = \frac{1}{T} \sum_{t=1}^T t(\mathbf{x})$. We have

$$\Pr \left\{ \Lambda_T \leq 4G_1 \sqrt{D_T \ln \frac{m}{\mu}} + 2G_1 D \ln \frac{m}{\mu} \right. \\ \left. + \Pr \left\{ D_T \leq \frac{D^2}{T} \right\} \geq 1 - \epsilon \right\},$$

where $m = \lceil 2 \log_2 T \rceil$ and $t(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T (\nabla f(\mathbf{x}_t) - \mathbf{g}(\mathbf{x}_t))^\top (\mathbf{x} - \mathbf{x}_t)$.

Proof of Theorem 2 The proof below follows from techniques used in Lemma 2 and Theorem 1. Since $F(\mathbf{x})$ is μ -strongly convex, we have

$$F(\mathbf{x}_t) - F(\mathbf{x}) \leq (\mathbf{x}_t - \mathbf{x})^\top \nabla F(\mathbf{x}_t) - \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2.$$

Combining the above inequality with the inequality in (8) and taking summation over all $t = 1, \dots, T$, we have

$$\frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{x})) \leq \frac{\frac{1}{2} \sum_{t=1}^T \|\mathbf{x}_t - \mathbf{x}\|_2^2 + T(G_1^2 + \mu G^2)}{BT} \\ + \frac{1}{T} \sum_{t=1}^T t(\mathbf{x}) - \frac{1}{2} D_T. \quad (23)$$

We substitute the bound in Lemma 6 into the above inequality with $\mathbf{x} = \mathbf{x}_*$. We consider two cases. In the first

case, we assume $D_T \leq D^2/T$. As a result, we have

$$\frac{1}{T} \sum_{t=1}^T t(\mathbf{x}_*) = \frac{1}{T} \sum_{t=1}^T (\nabla f(\mathbf{x}_t) - \mathbf{g}(\mathbf{x}_t))^\top (\mathbf{x}_* - \mathbf{x}_t) \\ \leq 2G_1 \sqrt{D_T} \leq 2G_1 D,$$

which together with the inequality in (23) leads to the bound

$$\frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{x}_*)) \leq 2G_1 D + BT.$$

In the second case, we assume

$$\frac{1}{T} \sum_{t=1}^T t(\mathbf{x}_*) \leq 4G_1 \sqrt{D_T \ln \frac{m}{\mu}} + 4G_1 D \ln \frac{m}{\mu} \\ \leq \frac{1}{2} D_T + \frac{8G_1^2}{\mu} + 4G_1 D \ln \frac{m}{\mu},$$

where the last step uses the fact $2\sqrt{ab} \leq a^2 + b^2$. We thus have

$$\frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{x}_*)) \leq \frac{8G_1^2}{\mu} + 2G_1 D \ln \frac{m}{\mu} + BT$$

Combing the results of the two cases, we have, with a probability $1 - \epsilon$,

$$\frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{x}_*)) \leq \frac{8G_1^2}{\mu} + 2G_1 D \ln \frac{m}{\mu} \\ + 2G_1 D + BT,$$

where $C = \frac{8G_1^2}{\mu} + 2G_1 D \ln \frac{m}{\mu} + 2G_1 D$. Following the same analysis, we have

$$f(\mathbf{x}_T) - f(\mathbf{x}_*) \leq \frac{\mu C}{T} + \frac{\mu \|\mathbf{x}_1 - \mathbf{x}_*\|_2^2}{2T} + \mu G^2$$

Let $\Delta_k = f(\mathbf{x}_k^1) - f(\mathbf{x}_*)$. By induction, we have

$$\Delta_{k+1} \leq \frac{\mu C}{T_k} + \frac{\mu \Delta_k}{2 T_k} + \mu G^2$$

Assume $\Delta_k \leq V_k$ and $\frac{\mu^2 G^2}{2^{k-2}}$, by plugging the values of k, T_k , we have

$$\Delta_{k+1} \leq \frac{V_k}{6} + \frac{V_k}{6} + \frac{V_k}{6} = \frac{V_k}{2} = V_{k+1}$$

where we use $T_1 \geq \max\{\frac{3C}{\mu G^2}, 9\}$ and $T_k \geq \max\{\frac{6\mu C}{V_k}, \frac{18\mu^2 G^2}{V_k}\}$ and $k = \frac{V_k}{6\mu G^2} = \frac{2\mu}{2^k(3)}$. This completes the proof of this theorem.

Proof of Lemma 3

To prove Lemma 3, we derive an inequality similar to Eq. (8); the rest proof of Lemma 3 is similar to that of Lemma 2.

Corollary 1. *Given a μ -strongly convex function $f(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$, and a sequence $\{\mathbf{x}_t\}$ defined by the update $\mathbf{x}_{t+1} = \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_t - \mathbf{g}(\mathbf{x}_t))\|_2^2 + g(\mathbf{x})$. Then for any \mathbf{x} , we have*

$$\begin{aligned} & \sum_{t=1}^T [f(\mathbf{x}_t) + g(\mathbf{x}_{t+1}) - f(\mathbf{x}) - g(\mathbf{x})] \\ \leq & \frac{\|\mathbf{x} - \mathbf{x}_1\|_2^2}{2} + \frac{1}{2} \sum_{t=1}^T \|\mathbf{g}(\mathbf{x}_t)\|_2^2 + \sum_{t=1}^T (\mathbf{x} - \mathbf{x}_t)^\top (\mathbf{g}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)) - \frac{1}{2} \sum_{t=1}^T \|\mathbf{x} - \mathbf{x}_{t+1}\|_2^2. \end{aligned}$$

Corollary 1 can be proved using techniques similar to the ones in [9] but with extra care on the stochastic gradient. As a consequence we have

$$\begin{aligned} & \frac{1}{T} \mathbb{E} \sum_{t=1}^T \hat{f}(\mathbf{x}_t) - \hat{f}(\mathbf{x}) \\ \leq & \frac{\mathbb{E}[\|\mathbf{x} - \mathbf{x}_1\|_2^2]}{2T} + (G_1^2 + G_2^2) + \frac{g(\mathbf{x}_1) - g(\mathbf{x}_{T+1})}{T} \end{aligned}$$

Proof of Lemma 4

The lemma is a corollary of results in [6] for general convex optimization. In particular, if we consider the stochastic composite optimization

$$F(\mathbf{x}) = \psi(\mathbf{x}) + g(\mathbf{x})$$

where $g(\mathbf{x})$ is a simple function such that its proximal mapping can be easily solved and $\psi(\mathbf{x})$ is only accessible through a stochastic oracle that returns a stochastic subgradient $\mathbf{g}(\mathbf{x})$. To state the convergence of ORDA for general convex problems, [6] makes the following assumptions: (i) $\mathbb{E}[\|\mathbf{g}(\mathbf{x}) - \mathbb{E}\mathbf{g}(\mathbf{x})\|_2^2] \leq M^2$ and (ii)

$$\|\mathbf{y}\|_2 - \|\mathbf{x}\|_2 - (\mathbf{y} - \mathbf{x})^\top \mathbf{c}(\mathbf{x}) \leq M \|\mathbf{y} - \mathbf{x}\|_2$$

When $\|\mathbf{c}(\mathbf{x})\|_2 \leq G$, the first inequality holds with $M = G$ and the second inequality holds with $M = 2G$. Applying to the augmented objective

$$F(\mathbf{x}) = f(\mathbf{x}) + [c(\mathbf{x})]_+ + g(\mathbf{x})$$

We note that $G = G_1$ and $M = 2(G_1 + G_2)$. Follow the inequality (26) in the appendix of [6], we obtain that

$$\mathbb{E}[F(\mathbf{x}_{T+2}) - F(\mathbf{x}_*)] \leq \frac{4\|\mathbf{x}_1 - \mathbf{x}_*\|_2^2}{\sqrt{T}} + \frac{2(G_1 + M)^2}{\sqrt{T}}$$

by using the Euclidean distance $V(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ and their notation $\mu = 1/c$, and noting that c is the inverse of their notation c . Then the second inequality in Lemma 4 can be proved similarly as for Lemma 2.

Proof of Theorem 3

Proof. Recall $\mu = 1/c$ ($c = G_1/c$) and $G = 3G_1 + 2G_2$. Let $V_k = \mu^2 G^2 / 2^{k-2}$. By the values of V_k and T_k we have

$$T_k = 2^{k+3} = \frac{32\mu^2 G^2}{V_k}, \quad k = \frac{\mu}{2^{(k-1)/2}} = \frac{V_k \sqrt{T_k}}{8\mu G^2}.$$

Define $\Delta_k = \hat{f}(\mathbf{x}_1^k) - \hat{f}(\mathbf{x}_*)$. We first prove the inequality

$$\mathbb{E}[\Delta_k] \leq V_k$$

by induction. It is true for $k = 1$ because of Lemma 5, $\mu > 1$ and $G^2 > G_1^2$. Now assume it is true for k and we prove it for $k+1$. For a random variable X measurable with respect to the randomness up to epoch $k+1$. Let $\mathbb{E}_k[X]$ denote the expectation conditioned on all the randomness up to epoch k . Following Lemma 2, we have

$$\mathbb{E}_k[\Delta_{k+1}] \leq \mu \frac{2G^2}{\sqrt{T_k}} + \frac{\mathbb{E}[4\|\mathbf{x}_1^k - \mathbf{x}_*\|_2^2]}{k\sqrt{T_k}} \quad (24)$$

Since $\Delta_k = f(\mathbf{x}_1^k) - f(\mathbf{x}_*) \geq \|\mathbf{x}_1^k - \mathbf{x}_*\|_2^2/2$ by the strong convexity, we have

$$\begin{aligned} \mathbb{E}[\Delta_{k+1}] & \leq \mu \frac{2G^2}{\sqrt{T_k}} + \frac{\mathbb{E}[8\Delta_k]}{k\sqrt{T_k}} \\ & = \frac{2\mu G^2}{\sqrt{T_k}} + \frac{V_k \mu}{k\sqrt{T_k}} = \frac{V_k}{4} + \frac{V_k}{4} = \frac{V_k}{2} \end{aligned}$$

where we use the fact $\mu/\sqrt{T_k} = V_k/(8\mu G^2)$ and $T_k = 32\mu^2 G^2/(V_k)$. Thus, we get

$$\mathbb{E}[f(\mathbf{x}_1^{k+1})] - f(\mathbf{x}_*) = \mathbb{E}[\Delta_{k+1}] \leq V_{k+1} = \frac{\mu^2 G^2}{2^{k+1}}$$

Note that the total number of epochs satisfies

$$\sum_{k=1}^{k^\dagger} (T_k + 1) = 16(2^{k^\dagger} - 1) + k^\dagger \leq T$$

By some reformulations, we complete the proof of this theorem. \square

Proof of Lemma 6

The proof of Lemma 6 is based on the *Bernstein Inequality for Martingales* [4]. We present its main result below for completeness.

Theorem 4. [Bernstein Inequality for Martingales] Let X_1, \dots, X_n be a bounded martingale difference sequence with respect to the filtration $\mathcal{F} = (\mathcal{F}_i)_{1 \leq i \leq n}$ and with $\|X_i\| \leq K$. Let

$$S_i = \sum_{j=1}^i X_j$$

be the associated martingale. Denote the sum of the conditional variances by

$$\Sigma_n^2 = \sum_{t=1}^n \mathbb{E} X_t^2 | \mathcal{F}_{t-1} ,$$

Then for all constants $t, \delta > 0$,

$$\Pr \max_{i=1}^n$$