

## A Proof of Lemma 4.1

*Proof.* By the definition of  $\mathbf{v}_B$ , we directly have

$$\begin{aligned}\mathbb{E}_B \mathbf{v}_B &= \mathbb{E}_B \left( \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla f_i(\tilde{\boldsymbol{\theta}}) + \nabla \mathcal{F}(\tilde{\boldsymbol{\theta}}) \right) \\ &= \frac{1}{|\mathcal{B}|} \mathbb{E}_B \sum_{i \in \mathcal{B}} \nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \frac{1}{|\mathcal{B}|} \mathbb{E}_B \sum_{i \in \mathcal{B}} \nabla f_i(\tilde{\boldsymbol{\theta}}) + \nabla \mathcal{F}(\tilde{\boldsymbol{\theta}}) \\ &= \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}) - \nabla \mathcal{F}(\tilde{\boldsymbol{\theta}}) + \nabla \mathcal{F}(\tilde{\boldsymbol{\theta}}) = \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}).\end{aligned}$$

Thus  $\mathbf{v}_B$  is an unbiased estimator of  $\nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})$ . Let  $i$  be an index sampled from  $\{1, \dots, n\}$  with equal probability, we define  $\mathbf{v}_i = \nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\tilde{\boldsymbol{\theta}}) + \tilde{\boldsymbol{\mu}}$ . Since all indices in  $\mathcal{B}$  are independently sampled from  $\{1, \dots, n\}$  with equal probability, we have

$$\mathbb{E}_B \|\mathbf{v}_B - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2 = \frac{1}{|\mathcal{B}|} \mathbb{E}_i \|\mathbf{v}_i - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2. \quad (\text{A.1})$$

We then proceed to bound R.H.S. of (A.1) from above as follows,

$$\begin{aligned}\mathbb{E}_i \|\mathbf{v}_i - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2 &\leq \mathbb{E}_i \|\nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\tilde{\boldsymbol{\theta}}) + \nabla \mathcal{F}(\tilde{\boldsymbol{\theta}}) - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2 \\ &\stackrel{(i)}{\leq} \mathbb{E}_i \|\nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\tilde{\boldsymbol{\theta}})\|^2 \\ &\stackrel{(ii)}{\leq} 2\mathbb{E}_i \|\nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2 + 2\mathbb{E}_i \|\nabla f_i(\tilde{\boldsymbol{\theta}}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2 \\ &= \frac{2}{n} \sum_{i=1}^n \|\nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2 + \frac{2}{n} \sum_{i=1}^n \|\nabla f_i(\tilde{\boldsymbol{\theta}}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2, \quad (\text{A.2})\end{aligned}$$

where (i) comes from  $\mathbb{E} \|\mathbf{u} - \mathbb{E} \mathbf{u}\|^2 = \mathbb{E} \|\mathbf{u}\|^2 - \|\mathbb{E} \mathbf{u}\|^2 \leq \mathbb{E} \|\mathbf{u}\|^2$ , and (2) comes from  $\|\mathbf{u} - \mathbf{w}\|^2 \leq 2\|\mathbf{u}\|^2 + 2\|\mathbf{w}\|^2$  for any random vectors  $\mathbf{u}$  and  $\mathbf{w}$ .

We then define  $g_i(\boldsymbol{\theta}) = f_i(\boldsymbol{\theta}) - f_i(\hat{\boldsymbol{\theta}}) - \nabla f_i(\hat{\boldsymbol{\theta}})^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$ . By Assumption 2.1, we have

$$\|\nabla g_i(\boldsymbol{\theta}') - \nabla g_i(\boldsymbol{\theta})\| = \|\nabla f_i(\boldsymbol{\theta}') - \nabla f_i(\boldsymbol{\theta})\| \leq T_{\max} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|, \quad (\text{A.3})$$

which implies that  $\nabla g_i(\boldsymbol{\theta})$  is also Lipschitz continuous. Thus we have

$$\begin{aligned}0 &= g_i(\hat{\boldsymbol{\theta}}) \leq \min_{\alpha} g_i(\boldsymbol{\theta} - \alpha \nabla g_i(\boldsymbol{\theta})) \\ &= \min_{\alpha} [g_i(\boldsymbol{\theta}) - \alpha \|\nabla g_i(\boldsymbol{\theta})\|^2 + \frac{T_{\max} \alpha^2}{2} \|\nabla g_i(\boldsymbol{\theta})\|^2], \quad (\text{A.4})\end{aligned}$$

where the first inequality comes from the fact that  $\hat{\boldsymbol{\theta}}$  is the minimizer to (A.3). Minimizing R.H.S. of (A.4), we have  $\alpha = 1/T_{\max}$ . Then (A.4) can be written as

$$0 \leq g_i(\boldsymbol{\theta}) - \frac{1}{2T_{\max}} \|\nabla g_i(\boldsymbol{\theta})\|^2. \quad (\text{A.5})$$

Combining the definition of  $g_i(\boldsymbol{\theta})$  and (A.5), we further have

$$\|\nabla f_i(\boldsymbol{\theta}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2 \leq 2T_{\max} \left[ f_i(\boldsymbol{\theta}) - f_i(\hat{\boldsymbol{\theta}}) - \nabla f_i(\hat{\boldsymbol{\theta}})^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right]. \quad (\text{A.6})$$

Taking the summation of (A.6) over  $i = 1, \dots, n$ , we obtain

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(\boldsymbol{\theta}) - \nabla f_i(\hat{\boldsymbol{\theta}})\|^2 \leq 2T_{\max} \left[ \mathcal{F}(\boldsymbol{\theta}) - \mathcal{F}(\hat{\boldsymbol{\theta}}) - \nabla \mathcal{F}(\hat{\boldsymbol{\theta}})^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right]. \quad (\text{A.7})$$

The optimality condition of  $\hat{\boldsymbol{\theta}}$  implies that there exists some  $\boldsymbol{\xi} \in \partial \mathcal{R}(\hat{\boldsymbol{\theta}})$  such that

$$\nabla \mathcal{F}(\hat{\boldsymbol{\theta}}) + \boldsymbol{\xi} = \mathbf{0}. \quad (\text{A.8})$$

Combining (A.8) with the convexity of  $\mathcal{R}$ , we have

$$\nabla \mathcal{F}(\hat{\boldsymbol{\theta}})^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) = -\boldsymbol{\xi}^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \leq \mathcal{R}(\boldsymbol{\theta}) - \mathcal{R}(\hat{\boldsymbol{\theta}}). \quad (\text{A.9})$$

Combining (A.1), (A.7), (A.9), and (A.2), we eventually obtain

$$\mathbb{E}_B \|\mathbf{v}_B - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2 \leq \frac{4T_{\max}}{|\mathcal{B}|} \left[ \mathcal{P}(\boldsymbol{\theta}^{(t-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) + \mathcal{P}(\tilde{\boldsymbol{\theta}}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \right]. \quad (\text{A.10})$$

□

## B Proof of Theorem 4.2

Before we proceed with the proof, we first define

$$\mathcal{T}_\eta(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}'}{\operatorname{argmin}} \frac{1}{2\eta} \|\boldsymbol{\theta}' - \boldsymbol{\theta}\|^2 + \mathcal{R}(\boldsymbol{\theta}') = \underset{\boldsymbol{\theta}'}{\operatorname{argmin}} \frac{1}{2\eta} \sum_{j=1}^k \|\boldsymbol{\theta}'_{\mathcal{G}_j} - \boldsymbol{\theta}_{\mathcal{G}_j}\|^2 + \sum_{j=1}^k r_j(\boldsymbol{\theta}'_{\mathcal{G}_j}),$$

where the last equality comes from the assumption that  $\mathcal{R}(\boldsymbol{\theta}')$  is block separable. Then by Assumption 2.3, we have

$$\mathcal{T}_\eta(\boldsymbol{\theta}) = (\mathcal{T}_\eta^1(\boldsymbol{\theta}_{\mathcal{G}_1})^T, \dots, \mathcal{T}_\eta^k(\boldsymbol{\theta}_{\mathcal{G}_k})^T)^T.$$

For any vector  $\mathbf{v} \in \mathbb{R}^d$ , we define  $\mathbf{v}^{\mathcal{G}_j} = (\mathbf{0}^T, \dots, \mathbf{v}_{\mathcal{G}_j}^T, \dots, \mathbf{0}^T)^T$ . It is easy to verify  $\mathbf{v} = \sum_{j=1}^k \mathbf{v}^{\mathcal{G}_j}$ , and  $\mathbf{v}^{\mathcal{G}_j}$  and  $\mathbf{v}^{\mathcal{G}_{j'}}$  are orthogonal to each other, for any  $j \neq j'$ , i.e.,  $(\mathbf{v}^{\mathcal{G}_j})^T \mathbf{v}^{\mathcal{G}_{j'}} = 0$ . We then define  $\bar{\boldsymbol{\theta}} = \mathcal{T}_\eta(\boldsymbol{\theta} - \eta \mathbf{v})$  and

$$\bar{\boldsymbol{\theta}}^{\mathcal{G}_j} = (\boldsymbol{\theta}_{\mathcal{G}_1}^T, \dots, [\mathcal{T}_\eta^j(\boldsymbol{\theta}_{\mathcal{G}_j} - \eta \mathbf{v}_{\mathcal{G}_j})]^T, \dots, \boldsymbol{\theta}_{\mathcal{G}_k}^T)^T.$$

We then introduce the following lemma:

**Lemma B.1.** Define  $\boldsymbol{\delta} = (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta})/\eta$  and  $\boldsymbol{\delta}^{\mathcal{G}_j} = (\bar{\boldsymbol{\theta}}^{\mathcal{G}_j} - \boldsymbol{\theta})/\eta$ . If the block index  $j$  is randomly selected from  $\{1, \dots, k\}$  with equal probability, then we have

$$\mathbb{E}_j[\boldsymbol{\delta}^{\mathcal{G}_j}] = \boldsymbol{\delta}/k \quad \text{and} \quad \mathbb{E}_j[\|\boldsymbol{\delta}^{\mathcal{G}_j}\|^2] = \|\boldsymbol{\delta}\|^2/k.$$

Moreover, taking  $\eta \leq 1/L_{\max}$ , we have

$$\begin{aligned} \mathbb{E}_j \left[ (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T \boldsymbol{\delta}^{\mathcal{G}_j} + \frac{\eta}{2} \|\boldsymbol{\delta}^{\mathcal{G}_j}\|^2 \right] \\ \leq \frac{1}{k} \mathcal{P}(\hat{\boldsymbol{\theta}}) + \frac{k-1}{k} \mathcal{P}(\boldsymbol{\theta}) - \mathbb{E}_j[\mathcal{P}(\bar{\boldsymbol{\theta}}^{\mathcal{G}_j})] + \frac{1}{k} (\mathbf{v} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T (\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}). \end{aligned}$$

The proof of Lemma B.1 is presented in Appendix D.

Now we proceed with the proof of Theorem 4.2. At the  $t$ -th iteration of the inner loop, we randomly sample a mini-batch  $\mathcal{B}$  and a block of coordinates  $\mathcal{G}_j$ . Define  $\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j} = (\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^{(t-1)})/\eta$ . We then have

$$\begin{aligned} \mathbb{E}_{\mathcal{B},j} \|\boldsymbol{\theta}^{(t)} - \hat{\boldsymbol{\theta}}\|^2 &= \mathbb{E}_{\mathcal{B},j} \|\boldsymbol{\theta}^{(t-1)} + \eta \boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j} - \hat{\boldsymbol{\theta}}\|^2 \\ &= \|\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}}\|^2 + 2\eta \mathbb{E}_{\mathcal{B},j} [(\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}})^T \boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}] + \eta^2 \mathbb{E}_{\mathcal{B},j} \|\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}\|^2 \\ &= \|\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}}\|^2 + 2\eta (\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}})^T \mathbb{E}_{\mathcal{B},j} [\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}] + \eta^2 \mathbb{E}_{\mathcal{B},j} \|\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}\|^2 \\ &= \|\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}}\|^2 + \mathbb{E}_{\mathcal{B}} \left[ 2\eta (\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}})^T \mathbb{E}_j [\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}] + \eta^2 \mathbb{E}_j \|\boldsymbol{\delta}_{\mathcal{B}}^{\mathcal{G}_j}\|^2 \right]. \quad (\text{B.1}) \end{aligned}$$

Define  $\bar{\boldsymbol{\theta}}_{\mathcal{B}} = \mathcal{T}_\eta(\boldsymbol{\theta}^{(t-1)} - \eta \mathbf{v}_{\mathcal{B}})$  and  $\mathbf{v}_{\mathcal{B}} = \sum_{i \in \mathcal{B}} \nabla f_i(\boldsymbol{\theta}^{(t-1)}) - \nabla f_i(\tilde{\boldsymbol{\theta}}) + \tilde{\boldsymbol{\mu}}$ . Then by applying Lemma B.1 to (B.1), we further have

$$\begin{aligned} \mathbb{E}_{\mathcal{B},j} \|\boldsymbol{\theta}^{(t)} - \hat{\boldsymbol{\theta}}\|^2 &= \|\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}}\|^2 \\ &\leq \frac{2\eta}{k} \mathbb{E}_{\mathcal{B}} [(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T (\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}})] + 2\eta \mathbb{E}_{\mathcal{B}} \left[ \frac{1}{k} \mathcal{P}(\hat{\boldsymbol{\theta}}) + \frac{k-1}{k} \mathcal{P}(\boldsymbol{\theta}) - \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) \right] \\ &\leq \frac{2\eta}{k} \mathbb{E}_{\mathcal{B}} [(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T (\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}})] \\ &\quad + 2\eta \mathbb{E}_{\mathcal{B}} \left[ \frac{1}{k} \mathcal{P}(\hat{\boldsymbol{\theta}}) + \frac{k-1}{k} \mathcal{P}(\hat{\boldsymbol{\theta}}) - \frac{k-1}{k} \mathcal{P}(\hat{\boldsymbol{\theta}}) + \frac{k-1}{k} \mathcal{P}(\boldsymbol{\theta}) - \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) \right] \\ &\leq \frac{2\eta}{k} \mathbb{E}_{\mathcal{B}} [(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T (\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}})] \\ &\quad - 2\eta \mathbb{E}_{\mathcal{B}} \left[ \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) - \frac{k-1}{k} (\mathcal{P}(\hat{\boldsymbol{\theta}}) - \mathcal{P}(\boldsymbol{\theta})) \right]. \quad (\text{B.2}) \end{aligned}$$

To bound  $\mathbb{E}_{\mathcal{B}}[(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}})]$ , we define  $\bar{\boldsymbol{\theta}} = \mathcal{T}_{\eta}(\boldsymbol{\theta}^{(t-1)} - \eta \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))$ . Then we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{B}}(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}}) \\ &= \mathbb{E}_{\mathcal{B}} \left[ (\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}} + \bar{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}}) \right] \\ &= \mathbb{E}_{\mathcal{B}} \left[ (\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}) + (\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T(\bar{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}}) \right] \\ &= \mathbb{E}_{\mathcal{B}} \left[ (\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))^T(\bar{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}}) \right], \end{aligned}$$

where the last equality comes from the fact  $\mathbb{E}_{\mathcal{B}}[\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})] = \mathbf{0}$ . By Cauchy-Schwarz inequality, we further have

$$\begin{aligned} & \mathbb{E}_{\mathcal{B}}[(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}})] \\ & \leq \mathbb{E}_{\mathcal{B}}\|\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\| \cdot \|\mathcal{T}_{\eta}(\boldsymbol{\theta} - \eta \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})) - \mathcal{T}_{\eta}(\boldsymbol{\theta} - \eta \mathbf{v}_{\mathcal{B}})\| \\ & \leq \mathbb{E}_{\mathcal{B}}\|\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\| \cdot \|(\boldsymbol{\theta} - \eta \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})) - (\boldsymbol{\theta} - \eta \mathbf{v}_{\mathcal{B}})\| \\ & = \eta \mathbb{E}_{\mathcal{B}}\|\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}^{(t-1)})\|^2, \end{aligned} \tag{B.3}$$

where the second inequality comes from the non-expansiveness of the proximal operator [11]. Combining (B.2) and (B.3), we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{B},j}\|\boldsymbol{\theta}^{(t)} - \hat{\boldsymbol{\theta}}\|^2 - \|\boldsymbol{\theta}^{(t-1)} - \hat{\boldsymbol{\theta}}\|^2 \\ & \leq -2\eta \mathbb{E}_{\mathcal{B}} \left[ \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) - \frac{k-1}{k} (\mathcal{P}(\hat{\boldsymbol{\theta}}) - \mathcal{P}(\boldsymbol{\theta})) \right] + \frac{2\eta^2}{k} \mathbb{E}_{\mathcal{B}}(\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta}))^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}_{\mathcal{B}}) \\ & \leq -2\eta \mathbb{E}_{\mathcal{B}} \left[ \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) - \frac{k-1}{k} (\mathcal{P}(\hat{\boldsymbol{\theta}}) - \mathcal{P}(\boldsymbol{\theta})) \right] + \frac{2\eta^2}{k} \mathbb{E}_{\mathcal{B}}\|\mathbf{v}_{\mathcal{B}} - \nabla \mathcal{F}(\boldsymbol{\theta})\|^2 \\ & \leq -2\eta \mathbb{E}_{\mathcal{B}} \left[ \mathbb{E}_j \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) - \frac{k-1}{k} (\mathcal{P}(\hat{\boldsymbol{\theta}}) - \mathcal{P}(\boldsymbol{\theta})) \right] \\ & \quad + \frac{8\eta^2 T_{max}}{k|\mathcal{B}|} (\mathcal{P}(\boldsymbol{\theta}^{(t-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) + \mathcal{P}(\tilde{\boldsymbol{\theta}}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \\ & \leq -2\eta \left[ \mathbb{E}_{\mathcal{B},j} \mathcal{P}(\bar{\boldsymbol{\theta}}_{\mathcal{B}}^{\mathcal{G}_j}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \right] + 2\eta \frac{k-1}{k} (\mathcal{P}(\boldsymbol{\theta}^{(t-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) \\ & \quad + \frac{8T_{max}\eta^2}{k|\mathcal{B}|} (\mathcal{P}(\boldsymbol{\theta}^{(k-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) + \mathcal{P}(\tilde{\boldsymbol{\theta}}) - \mathcal{P}(\hat{\boldsymbol{\theta}}), \end{aligned} \tag{B.4}$$

where the third inequality comes from Lemma 4.1.

At the  $s$ -th iteration of the outer loop, we have

$$\boldsymbol{\theta}^{(0)} = \tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}^{(s-1)} \quad \text{and} \quad \tilde{\boldsymbol{\theta}}^{(s)} = \frac{1}{m} \sum_{t=1}^m \boldsymbol{\theta}^{(t)}. \tag{B.5}$$

Thus summing (B.4) over  $t = 1, \dots, m$  and taking expectation with respect to  $\mathcal{B}$  and  $j$  over all iterations of the inner loop, we obtain

$$\begin{aligned} & \mathbb{E}\|\boldsymbol{\theta}^{(m)} - \hat{\boldsymbol{\theta}}\|^2 - \|\boldsymbol{\theta}^{(0)} - \hat{\boldsymbol{\theta}}\|^2 + 2\eta \sum_{t=1}^m (\mathbb{E} \mathcal{P}(\boldsymbol{\theta}^{(t)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) \\ & \leq \frac{8T_{max}\eta^2/|\mathcal{B}| + 2\eta(k-1)}{k} \sum_{t=1}^{m-1} (\mathbb{E} \mathcal{P}(\boldsymbol{\theta}^{(t)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) + \frac{8T_{max}\eta^2(m+1)}{k|\mathcal{B}|} (\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) \\ & \leq \frac{8T_{max}\eta^2/|\mathcal{B}| + 2\eta(k-1)}{k} \sum_{t=1}^m (\mathbb{E} \mathcal{P}(\boldsymbol{\theta}^{(t)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) \\ & \quad + \frac{8T_{max}\eta^2(m+1)}{k|\mathcal{B}|} (\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})), \end{aligned} \tag{B.6}$$

where the last inequality comes from  $\mathbb{E}\mathcal{P}(\boldsymbol{\theta}^{(t)}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \geq 0$ . By rearranging (B.6), we obtain

$$\begin{aligned} & 2\eta \left( \frac{1 - 4\eta T_{\max}/|\mathcal{B}|}{k} \right) \sum_{t=1}^m [\mathbb{E}\mathcal{P}(\boldsymbol{\theta}^{(t)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})] \\ & \leq \|\boldsymbol{\theta}^{(0)} - \hat{\boldsymbol{\theta}}\|^2 + \frac{8T_{\max}\eta^2(m+1)}{k|\mathcal{B}|} \left( \mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \right) \\ & \leq \frac{2}{\mu} (\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})) + \frac{8T_{\max}\eta^2(m+1)}{k|\mathcal{B}|} \left( \mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \right), \end{aligned} \quad (\text{B.7})$$

where the last inequality comes from the strong convexity of  $\mathcal{P}$  and  $\tilde{\boldsymbol{\theta}}^{(s-1)} = \boldsymbol{\theta}^{(0)}$ . By the convexity of  $\mathcal{P}$  again, we have

$$\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s)}) \leq \frac{1}{m} \sum_{t=1}^m \mathcal{P}(\boldsymbol{\theta}^{(t)}). \quad (\text{B.8})$$

Therefore combining (B.7) and (B.8), we obtain

$$\begin{aligned} & 2\eta \left( \frac{1 - 4\eta T_{\max}/|\mathcal{B}|}{k} \right) m [\mathbb{E}\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})] \\ & \leq \left( \frac{2}{\mu} + \frac{8T_{\max}\eta^2(m+1)}{k|\mathcal{B}|} \right) [\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})]. \end{aligned} \quad (\text{B.9})$$

Define

$$\alpha = \left( \frac{k}{\mu\eta(1 - 4\eta T_{\max}/|\mathcal{B}|)m} + \frac{4\eta T_{\max}/|\mathcal{B}|(m+1)}{(1 - 4\eta T_{\max}/|\mathcal{B}|)m} \right).$$

Then (B.9) implies

$$\mathbb{E}\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s)}) - \mathcal{P}(\hat{\boldsymbol{\theta}}) \leq \alpha [\mathcal{P}(\tilde{\boldsymbol{\theta}}^{(s-1)}) - \mathcal{P}(\hat{\boldsymbol{\theta}})]. \quad (\text{B.10})$$

By applying (B.10) recursively and setting  $|\mathcal{B}| = T_{\max}/L$ , we complete the proof.

## C Proof of Corollary 4.3

*Proof.* Note that choosing the suggested  $\eta$ ,  $m$ , and  $\mathcal{B}$  guarantees

$$\alpha = \frac{k}{\mu\eta(1 - 4\eta L_{\max})m} + \frac{4\eta L_{\max}(m+1)}{(1 - 4\eta L_{\max})m} \leq 2/3.$$

By Markov inequality, we have

$$\mathbb{P} \left( P(\tilde{\boldsymbol{\theta}}^{(s)}) - P(\hat{\boldsymbol{\theta}}) \geq \epsilon \right) \stackrel{(i)}{\leq} \frac{1}{\epsilon} \left( \mathbb{E}P(\tilde{\boldsymbol{\theta}}^{(s)}) - P(\hat{\boldsymbol{\theta}}) \right) \stackrel{(ii)}{\leq} \frac{1}{\epsilon} (2/3)^s \left( P(\tilde{\boldsymbol{\theta}}^{(0)}) - P(\boldsymbol{\theta}_*) \right) \leq \rho.$$

where (i) comes from Theorem 4.2, and (ii) comes from choosing the suggested  $s$ .  $\square$

## D Proof of Lemma B.1

*Proof.* Since  $\boldsymbol{\delta} = \sum_{j=1}^k \boldsymbol{\delta}^{\mathcal{G}_j}$ , we directly have

$$\mathbb{E}_j[\boldsymbol{\delta}^{\mathcal{G}_j}] = \boldsymbol{\delta}/k. \quad (\text{D.1})$$

Since  $\boldsymbol{v}^{\mathcal{G}_j}$  and  $\boldsymbol{v}^{\mathcal{G}_{j'}}$  are orthogonal to each other for any  $j \neq j'$ , we have  $\|\boldsymbol{\delta}\|^2 = \sum_{j=1}^k \|\boldsymbol{\delta}^{\mathcal{G}_j}\|^2$ . Taking expectation over the block index  $j$ , we directly have

$$\mathbb{E}_j[\|\boldsymbol{\delta}^{\mathcal{G}_j}\|^2] = \|\boldsymbol{\delta}\|^2/k. \quad (\text{D.2})$$

Since  $\mathcal{F}$  and  $\mathcal{R}$  are convex, we have

$$\mathcal{R}(\hat{\boldsymbol{\theta}}) \geq \mathcal{R}(\bar{\boldsymbol{\theta}}) + \boldsymbol{\xi}^T(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}}), \quad (\text{D.3})$$

where  $\xi \in \partial \mathcal{R}(\bar{\theta})$ , and

$$\mathcal{F}(\hat{\theta}) \geq \mathcal{F}(\theta) + \nabla \mathcal{F}(\theta)^T (\hat{\theta} - \theta). \quad (\text{D.4})$$

Combining (D.3) and (D.4), we have

$$\mathcal{P}(\hat{\theta}) = \mathcal{F}(\hat{\theta}) + \mathcal{R}(\hat{\theta}) \geq \mathcal{F}(\theta) + \nabla \mathcal{F}(\theta)^T (\hat{\theta} - \theta) + R(\bar{\theta}) + \xi^T (\hat{\theta} - \bar{\theta}). \quad (\text{D.5})$$

Since  $\mathcal{R}$  is block separable, we have

$$\begin{aligned} \mathbb{E}_j \mathcal{R}(\bar{\theta}^{\mathcal{G}_j}) &= \frac{1}{k} \sum_{j=1}^k \mathcal{R}(\theta + \eta \delta^{\mathcal{G}_j}) \\ &= \frac{1}{k} \sum_{j=1}^k \left( \sum_{l \neq j} r_l(\theta_{\mathcal{G}_l}) + r_j(\theta_{\mathcal{G}_j} + \eta \delta_{\mathcal{G}_j}) \right) \\ &= \frac{1}{k} \left[ (k-1) \sum_{j=1}^k r_j(\theta_{\mathcal{G}_j}) + \sum_{j=1}^k r_j(\theta_{\mathcal{G}_j} + \eta \delta_{\mathcal{G}_j}) \right] \\ &= \frac{k-1}{k} \mathcal{R}(\theta) + \frac{1}{k} \mathcal{R}(\bar{\theta}). \end{aligned} \quad (\text{D.6})$$

By Assumption 2.1, we have

$$\begin{aligned} \mathbb{E}_j \mathcal{F}(\bar{\theta}^{\mathcal{G}_j}) &\stackrel{(i)}{\leq} \mathcal{F}(\theta) + \mathbb{E}_j \left( \nabla \mathcal{F}(\theta)^T (\bar{\theta}^{\mathcal{G}_j} - \theta) + \frac{L_{max}}{2} \|\bar{\theta}^{\mathcal{G}_j} - \theta\|^2 \right) \\ &\stackrel{(ii)}{=} \mathcal{F}(\theta) + \mathbb{E}_j \left( \eta \nabla \mathcal{F}(\theta)^T \delta^{\mathcal{G}_j} + \frac{\eta^2 L_{max}}{2} \|\delta^{\mathcal{G}_j}\|^2 \right) \\ &\stackrel{(iii)}{=} \mathcal{F}(\theta) + \frac{1}{k} \left( \eta \nabla \mathcal{F}(\theta)^T \delta + \frac{\eta^2 L_{max}}{2} \|\delta\|^2 \right) \\ &= \frac{k-1}{k} \mathcal{F}(\theta) + \frac{1}{k} \left( \mathcal{F}(\theta) + \eta \nabla \mathcal{F}(\theta)^T \delta + \frac{\eta^2 L_{max}}{2} \|\delta\|^2 \right), \end{aligned} \quad (\text{D.7})$$

where (i) comes from the fact that  $\bar{\theta}^{\mathcal{G}_j}$  and  $\theta$  are identical except the  $j$ -th block of coordinates, (ii) comes from the definition of  $\delta^{\mathcal{G}_j}$ , and (iii) comes from (D.1) and (D.2).

By rearranging (D.7), we have

$$\mathcal{F}(\theta) \geq k \mathbb{E}_j \mathcal{F}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{F}(\theta) - \nabla \mathcal{F}(\theta)^T (\bar{\theta} - \theta) - \frac{\eta^2 L_{max}}{2} \|\delta\|^2. \quad (\text{D.8})$$

Combining (D.5), (D.6), and (D.8), we further have

$$\begin{aligned} \mathcal{P}(\hat{\theta}) &\stackrel{(i)}{\geq} k \mathbb{E}_j \mathcal{F}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{F}(\theta) - \nabla \mathcal{F}(\theta)^T (\bar{\theta} - \theta) - \frac{\eta^2 L_{max}}{2} \|\delta\|^2 \\ &\quad + \nabla \mathcal{F}(\theta)^T (\hat{\theta} - \theta) + \mathcal{R}(\bar{\theta}) + \xi(\hat{\theta} - \bar{\theta}) \\ &= k \mathbb{E}_j \mathcal{F}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{F}(\theta) + \nabla \mathcal{F}(\theta)^T (\hat{\theta} - \bar{\theta}) - \frac{\eta^2 L_{max}}{2} \|\delta\|^2 + \mathcal{R}(\bar{\theta}) + \xi(\hat{\theta} - \bar{\theta}) \\ &\stackrel{(ii)}{=} k \mathbb{E}_j \mathcal{F}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{F}(\theta) + \nabla \mathcal{F}(\theta)^T (\hat{\theta} - \bar{\theta}) - \frac{\eta^2 L_{max}}{2} \|\delta\|^2 \\ &\quad + k \mathbb{E}_j \mathcal{R}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{R}(\theta) + \xi(\hat{\theta} - \bar{\theta}) \\ &= k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) - \frac{\eta^2 L_{max}}{2} \|\delta\|^2 + (\nabla \mathcal{F}(\theta) + \xi)^T (\hat{\theta} - \bar{\theta}) \\ &\stackrel{(iii)}{\geq} k \mathbb{E}_j \mathcal{P}(\bar{\theta}_{ij}) - (k-1) \mathcal{P}(\theta) - \frac{\eta}{2} \|\delta\|^2 + (\nabla \mathcal{F}(\theta) + \xi)^T (\hat{\theta} - \bar{\theta}) \end{aligned} \quad (\text{D.9})$$

where (ii) comes from (D.8), (ii) comes from (D.6), and (iii) comes from  $\eta \leq 1/L_{max}$ .

By the definition of  $\bar{\theta}$ , we have

$$\bar{\theta} = \underset{\theta'}{\operatorname{argmin}} \frac{1}{2} \|\theta' - (\theta - \eta v)\|^2 + \eta \mathcal{R}(\theta'). \quad (\text{D.10})$$

The optimality condition of (D.10) implies that there exists some  $\xi \in \partial \mathcal{R}(\bar{\theta})$  satisfying

$$\bar{\theta} - (\theta - \eta v) + \eta \xi = \mathbf{0},$$

which implies  $\xi = -\delta - v$ . Then by (D.9), we have

$$\begin{aligned} \mathcal{P}(\hat{\theta}) &\geq k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) - \frac{\eta}{2} \|\delta\|^2 + (\nabla \mathcal{F}(\theta) + \xi)(\hat{\theta} - \bar{\theta}) \\ &= k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) - \frac{\eta}{2} \|\delta\|^2 + (\nabla \mathcal{F}(\theta) - \delta - v)^T (\hat{\theta} - \bar{\theta}) \\ &= k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) - \frac{\eta}{2} \|\delta\|^2 \\ &\quad - (v - \nabla \mathcal{F}(\theta))^T (\hat{\theta} - \bar{\theta}) - \delta^T (\hat{\theta} - \theta) - \delta^T (\theta - \bar{\theta}). \end{aligned} \quad (\text{D.11})$$

Since  $\theta - \bar{\theta} = \eta \delta$ , (D.11) further implies

$$\begin{aligned} \mathcal{P}(\hat{\theta}) &\leq k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) - \frac{\eta}{2} \|\delta\|^2 \\ &\quad - (v - \nabla \mathcal{F}(\theta))^T (\hat{\theta} - \bar{\theta}) - \delta^T (\hat{\theta} - \theta) + \eta \|\delta\|^2 \\ &= k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) + \frac{\eta}{2} \|\delta\|^2 - (v - \nabla \mathcal{F}(\theta))^T (\hat{\theta} - \bar{\theta}) - \delta^T (\hat{\theta} - \theta) \\ &= k \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) - (k-1) \mathcal{P}(\theta) + \frac{k\eta}{2} \mathbb{E}_j \|\delta^{\mathcal{G}_j}\|^2 \\ &\quad - (v - \nabla \mathcal{F}(\theta))^T (\hat{\theta} - \bar{\theta}) - k \mathbb{E}_j (\hat{\theta} - \theta)^T \delta^{\mathcal{G}_j}, \end{aligned} \quad (\text{D.12})$$

where the last equality comes from (D.1) and (D.2). By rearranging (D.12), we obtain

$$\begin{aligned} \mathbb{E}_j (\theta - \hat{\theta})^T \delta^{\mathcal{G}_j} + \frac{\eta}{2} \mathbb{E}_j \|\delta^{\mathcal{G}_j}\|^2 &\leq \frac{1}{k} \mathcal{P}(\hat{\theta}) - \mathbb{E}_j \mathcal{P}(\bar{\theta}^{\mathcal{G}_j}) \\ &\quad + \frac{k-1}{k} \mathcal{P}(\theta) + \frac{1}{k} (v - \nabla \mathcal{F}(\theta))^T (\hat{\theta} - \bar{\theta}), \end{aligned} \quad (\text{D.13})$$

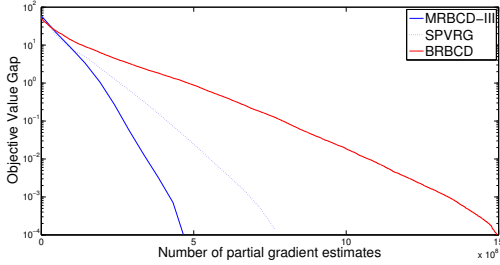
which completes the proof.  $\square$

## E Numerical Simulations

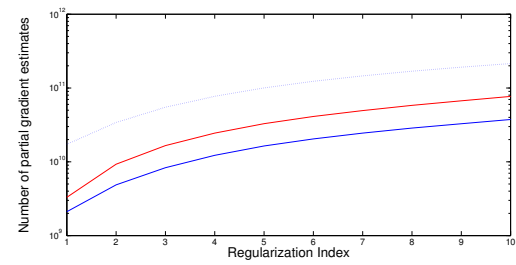
## F Real Data Experimental Results

Table E.1: Quantitive comparison of different methods on the simulated dataset for a sequence of regularization parameters. All three methods attains similar objective values for each regularization parameter, but MRBCD-III requires fewer partial gradient estimates than SPVRG and BRBCD.

|           |                |                |                |                |                |                |                |                |                |                |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| MRBCD     | $\lambda_1$    | $\lambda_2$    | $\lambda_3$    | $\lambda_4$    | $\lambda_5$    | $\lambda_6$    | $\lambda_7$    | $\lambda_8$    | $\lambda_9$    | $\lambda_{10}$ |
| # of P.G. | 29.32e5        | 61.12e5        | 93.81e5        | 126.42e5       | 159.1e5        | 192.2e5        | 225.0e5        | 260.1e5        | 300.6e5        | 343.2e5        |
| O.V.G.    | 9.23e-14       | 7.10e-14       | 7.45e-14       | 7.99e-14       | 7.81e-14       | 4.97e-14       | 4.61e-14       | 6.39e-14       | 4.26e-14       | 3.90e-14       |
| Reg.      | $\lambda_{11}$ | $\lambda_{12}$ | $\lambda_{13}$ | $\lambda_{14}$ | $\lambda_{15}$ | $\lambda_{16}$ | $\lambda_{17}$ | $\lambda_{18}$ | $\lambda_{19}$ | $\lambda_{20}$ |
| # of P.G. | 387.9e5        | 433.4e5        | 478.0e5        | 522.7e5        | 566.9e5        | 610.2e5        | 653.0e5        | 695.5e5        | 738.0e5        | 780.0e5        |
| O.V.G.    | 1.77e-14       | 2.48e-14       | 1.42e-14       | 3.55e-15       | 3.67e-15       | 4.67e-15       | 5.46e-15       | 5.57e-15       | 2.66e-15       | 1.78e-15       |
| SPVRG     | $\lambda_1$    | $\lambda_2$    | $\lambda_3$    | $\lambda_4$    | $\lambda_5$    | $\lambda_6$    | $\lambda_7$    | $\lambda_8$    | $\lambda_9$    | $\lambda_{10}$ |
| # of P.G. | 270.6e5        | 548.4e5        | 817.8e5        | 1074e5         | 1328e5         | 1586e5         | 1845e5         | 2133e5         | 2441e5         | 2776e5         |
| O.V.G.    | 8.57-14        | 9.43e-14       | 6.65e-14       | 9.12e-14       | 6.39e-14       | 4.97e-14       | 4.61e-14       | 6.39e-14       | 4.26e-14       | 0.461e-14      |
| Reg.      | $\lambda_{11}$ | $\lambda_{12}$ | $\lambda_{13}$ | $\lambda_{14}$ | $\lambda_{15}$ | $\lambda_{16}$ | $\lambda_{17}$ | $\lambda_{18}$ | $\lambda_{19}$ | $\lambda_{20}$ |
| # of P.G. | 3113e5         | 3454e5         | 3791e5         | 4124e5         | 4456e5         | 4782e5         | 5106e5         | 5425e5         | 5741e5         | 6053e5         |
| O.V.G.    | 3.90e-14       | 1.42e-14       | 3.19e-14       | 1.42e-14       | 8.88e-15       | 5.33e-15       | 3.55e-15       | 7.57e-15       | 4.44e-15       | 2.66e-15       |
| BRBCD     | $\lambda_1$    | $\lambda_2$    | $\lambda_3$    | $\lambda_4$    | $\lambda_5$    | $\lambda_6$    | $\lambda_7$    | $\lambda_8$    | $\lambda_9$    | $\lambda_{10}$ |
| # of P.G. | 43.50e5        | 95.80e5        | 153.2e5        | 209.5e5        | 264.8e5        | 320.4e5        | 375.7e5        | 435.4e5        | 508.8e5        | 585.2e5        |
| O.V.G.    | 5.68e-14       | 8.52e-14       | 7.81e-14       | 7.10e-14       | 7.10e-14       | 3.90e-14       | 4.26e-14       | 3.55e-14       | 5.68e-14       | 3.19e-14       |
| Reg.      | $\lambda_{11}$ | $\lambda_{12}$ | $\lambda_{13}$ | $\lambda_{14}$ | $\lambda_{15}$ | $\lambda_{16}$ | $\lambda_{17}$ | $\lambda_{18}$ | $\lambda_{19}$ | $\lambda_{20}$ |
| # of P.G. | 663.8e5        | 743.2e5        | 820.8e5        | 897.2e5        | 974.2e5        | 1050e5         | 1126e5         | 1201e5         | 1275e5         | 1356e5         |
| O.V.G.    | 7.11e-15       | 2.48e-14       | 1.42e-14       | 5.33e-15       | 3.55e-15       | 5.33e-15       | 3.55e-15       | 4.44e-15       | 1.78e-15       | 1.78e-15       |



(a) Comparison between different methods for a sin-



(b) Comparison between different methods for a sequence of regularization parameters.

Figure F.1: [a] The vertical axis corresponds to objective value gaps  $\mathcal{P}(\theta) - \mathcal{P}(\hat{\theta})$  in log scale. The horizontal axis corresponds to numbers of partial gradient estimates. [b] The horizontal axis corresponds to indices of regularization parameters. The vertical axis corresponds to numbers of partial gradient estimates in log scale. We see that MRBCD attains the best performance among all methods for both settings