

6 Appendix-A. Convergence Analysis

6.1 Convergence of Randomized Block Coordinate Descent

We first establish the linear convergence of Randomized Block Coordinate Descent (RBCD) when $L_n(\cdot)$ is smooth in the sense that its first derivative $L'_n(\cdot)$ is Lipschitz-continuous with parameter M_L , which then implies $L_n^*(\alpha_n)$ is strongly convex with parameter $1/M_L$.

Theorem 2-1 (Dual-RBCD for Smooth Loss). *Let the sequence $\{\alpha^s\}_{s=1}^\infty$ be the iterates produced by RBCD in the inner loop of Algorithm 2, and K be the number of blocks. Denote $\tilde{F}^*(\alpha)$ as the dual objective function of (18) and \tilde{F}_{opt}^* the optimal function value of (18). Then with probability $1 - \rho$,*

$$\tilde{F}^*(\alpha^s) - \tilde{F}_{opt}^* \leq \epsilon, \text{ for } s \geq \frac{K}{1 - c_1} \log\left(\frac{\tilde{F}^*(\alpha^0) - \tilde{F}_{opt}^*}{\rho\epsilon}\right) \quad (24)$$

if $L_n(\cdot)$ is smooth, where $0 < c_1 < 1$ is a constant depends on the smoothness parameter of $L_n(\cdot)$.

Proof. This is a special case of theorem 6 and theorem 4 in [13], where they consider composite objective function of the form

$$F(\alpha) = f(\alpha) + \Psi(\alpha), \quad (25)$$

where $f(\alpha)$ is a convex, smooth function, and $\Psi(\alpha)$ is a convex, block-separable function. In our case,

$$f(\alpha) = \tilde{R}^*\left(-\sum_{n=1}^N \Phi_n^T \alpha_n\right), \quad \Psi(\alpha) = \sum_{n=1}^N L_n^*(\alpha_n). \quad (26)$$

Note $\tilde{R}^*(\cdot)$ is smooth w.r.t. α_{B_k} with parameter $M_R = \eta_t B^2$, where $B \geq \|\Phi_{B_k}\|_2$ is an upper bound on the ℓ_2 -norm of each block's design matrix. If the loss function $L_n(\cdot)$ is smooth with parameter M_L , by Theorem 1, $\Psi(\alpha)$ is strongly convex with parameter $1/M_L$, and thus, based on [theorem 6, 21], (24) holds with

$$c_1 = \begin{cases} 1 - \frac{1}{4M_R M_L} & , \text{ if } M_R M_L \geq \frac{1}{2} \\ M_R M_L & , \text{ o.w..} \end{cases} \quad (27)$$

□

For some important classes of ERM, such as Support Vector Machine (SVM) and its variants (e.g. Multiclass, Structural SVM), $L_n(\alpha_n)$ is not smooth but piecewise-linear. In the following, we show that the linear convergence of RBCD holds for any loss $L_n(\alpha_n)$ with polyhedral epigraph if $R(\mathbf{w})$ is also polyhedral or smooth. The proof utilizes a restricted version of Strong Convexity called Constant Nullspace Strong Convexity [20, 22] and obtains a much tighter bound for RBCD than the bound proved in [20] for general feasible descent method. The proof follows is a generalization of that in [25] for proving linear convergence of RCD applied to the Augmented Lagrangian of Linear Program.

The augmented dual objective function (25), after some algebraic rearrangement, is equivalent to

$$\min_{\alpha, \mu} \sum_{n=1}^N L_n^*(\alpha_n) + R^*(-\mu) + \frac{\eta_t}{2} \left\| \sum_{n=1}^N \Phi_n^T \alpha_n - \mu + \mathbf{w}^t / \eta_t \right\|^2 \quad (28)$$

up to a constant. For $L_n^*(\alpha_n)$, $R^*(-\mu)$ being polyhedral, their epigraphs $\mathbf{epi}(L_n)$, $\mathbf{epi}(R)$ are polyhedrons and thus (28) can be also formulated as

$$\begin{aligned} \min_{\alpha, \mu, t, s} \quad & \sum_{n=1}^N t_n + r + \frac{\eta_t}{2} \left\| \sum_{n=1}^N \Phi_n^T \alpha_n - \mu + \mathbf{w}^t / \eta_t \right\|^2 \\ \text{s.t.} \quad & (\alpha_n, t_n) \in \mathbf{epi}(L_n) \\ & (\mu, r) \in \mathbf{epi}(R), \end{aligned} \quad (29)$$

which is of the form

$$\begin{aligned} \min_{\mathbf{x}} \quad & F(\mathbf{x}) = g(\bar{\Phi}^T \mathbf{x}) + \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & \mathbf{x} \in \mathcal{P} \end{aligned} \quad (30)$$

where $g(\mathbf{z}) = \frac{\eta_t}{2} \|\mathbf{z} + \mathbf{w}^t / \eta_t\|^2$ is a strongly convex function, \mathcal{P} is a polyhedral set

$$\mathcal{P} = \{(\boldsymbol{\alpha}, \mathbf{t}, \boldsymbol{\mu}, r) \mid (\boldsymbol{\alpha}_n, t_n) \in \text{epi}(L_n), (\boldsymbol{\mu}, r) \in \text{epi}(R)\},$$

and

$$\mathbf{x} = \begin{bmatrix} \boldsymbol{\alpha} \\ \mathbf{t} \\ \boldsymbol{\mu} \\ r \end{bmatrix} \quad \bar{\Phi} = \begin{bmatrix} \Phi \\ O_{N,d} \\ -I_{d,d} \\ O_{1,d} \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \\ \mathbf{0} \\ 1 \end{bmatrix}.$$

We will use I_α, I_t, I_μ and I_s denote the set of variable indexes j that correspond to $\boldsymbol{\alpha}, \mathbf{t}, \boldsymbol{\mu}$ and r respectively. For this type of objective function, we can show that the set of optimal solutions is a polyhedron defined by the following Lemma.

Lemma 1 (Lemma 4.2 of [20]). *The optimal solutions to problem (30) forms a polyhedral set*

$$\mathcal{S} = \{\mathbf{x} \mid \bar{\Phi}^T \mathbf{x} = \mathbf{p}^*, \mathbf{c}^T \mathbf{x} = q^*, \mathbf{x} \in \mathcal{P}\} \quad (31)$$

for some unique \mathbf{p}^*, q^* .

Furthermore, we can utilize the Hoffman's bound (defined in the following) to bound the distance of any point \mathbf{x} to the optimal solution set \mathcal{S} .

Lemma 2 (Hoffman's Bound). *Let $\mathcal{S} = \{\mathbf{x} \in \mathbb{R}^d \mid A\mathbf{x} \leq \mathbf{c}, E\mathbf{x} = \mathbf{c}\}$ be a polyhedral set. Then for any point $\mathbf{x} \in \mathbb{R}^d$,*

$$\|\mathbf{x} - \Pi_{\mathcal{S}}(\mathbf{x})\|_2^2 \leq \theta \left\| \begin{bmatrix} A\mathbf{x} - \mathbf{c} \\ E\mathbf{x} - \mathbf{c} \end{bmatrix}_+ \right\|_2^2 \quad (32)$$

where $\Pi_{\mathcal{S}}(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathcal{S}} \|\mathbf{y} - \mathbf{x}\|$ is the projection of \mathbf{x} to the set \mathcal{S} , and $\theta > 0$ is a constant depending on the polyhedral set \mathcal{S} .

Proof. The Hoffman's bound first appears in [4] and a proof for the ℓ_2 -norm's version (32) and the definition of the constant $\theta(\mathcal{S})$ can be found in [20] (lemma 4.3). \square

By Lemma 2, for any feasible $\mathbf{x} \in \mathcal{P}$, we obtain error bound

$$\|\mathbf{x} - \mathbf{x}^*\|^2 \leq \theta(\mathcal{S}) (\|\bar{\Phi}^T \mathbf{x} - \mathbf{p}^*\|^2 + \|\mathbf{c}^T \mathbf{x} - q^*\|^2), \quad (33)$$

which plays a crucial role in the proof of linear convergence.

The RBCD algorithm performed on (25) can be considered as minimizing (29) w.r.t. a block of dual variables $\{(\boldsymbol{\alpha}_n, t_n)\}_{n \in B_k}$ together with $(\boldsymbol{\mu}, s)$, while fixing all other variables $\{(\boldsymbol{\alpha}_n, t_n)\}_{n \notin B_k}$. In the following, we show that each block minimization step leads to a significant progress.

Lemma 3 (Descent Amount). *The expected descent amount for each Block Minimization step of Algorithm 2 has*

$$\mathbb{E}[F(\mathbf{x}^{k+1})] - F(\mathbf{x}^k) \leq \frac{1}{K} \left(\min_{\boldsymbol{\delta}} h(\mathbf{x}^k + \boldsymbol{\delta}) + \langle \nabla F(\mathbf{x}^k), \boldsymbol{\delta} \rangle + \frac{M\eta_t}{2} \|\boldsymbol{\delta}\|^2 \right), \quad (34)$$

where

$$h(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \mathcal{P} \\ \infty, & \text{o.w.} \end{cases} \quad (35)$$

and $M \geq \max_{k \in [K]} \|\Phi_{B_k}\|_2^2$ denotes a bound on the spectral norm of each block's design matrix.

Proof. First, notice that RBCD optimizes the function form of only variable $\boldsymbol{\alpha}$ while maintains other variables $(\mathbf{t}, \boldsymbol{\mu}, s)$ as their optimal values in each block minimization step, so we have

$$\mathbf{0} = \min_{\boldsymbol{\mu}, r} h(\mathbf{x}) + \nabla F(\mathbf{x}^s)^T (\mathbf{x} - \mathbf{x}^s) + \frac{M\eta_t}{2} \|\mathbf{x} - \mathbf{x}^s\|^2. \quad (36)$$

The algorithm picks coordinate uniformly from $\{(\alpha_{B_k}, t_{B_k})\}_{k=1}^K$ to update. Since the constant M upper bounds $\|\nabla_{\alpha_{B_k}, t_{B_k}} F(\mathbf{x})\|_2^2$, we have

$$\begin{aligned} F(\mathbf{x}^{s+1}) - F(\mathbf{x}^s) &= F(\boldsymbol{\alpha}^{s+1}, t^*(\boldsymbol{\alpha}^{s+1}), \boldsymbol{\mu}^*(\boldsymbol{\alpha}^{s+1}), r^*(\boldsymbol{\alpha}^{s+1})) - F(\mathbf{x}^s) \\ &\leq F(\boldsymbol{\alpha}^{s+1}, t^*(\boldsymbol{\alpha}^{s+1}), \boldsymbol{\mu}^s, r^s) - F(\mathbf{x}^s) \\ &\leq \min_{\boldsymbol{\delta}_{B_k}} h(\mathbf{x}^s + \boldsymbol{\delta}_{B_k}) + \nabla_{B_k} F(\mathbf{x}^k)^T \boldsymbol{\delta}_{B_k} + \frac{M\eta_t}{2} \|\boldsymbol{\delta}_{B_k}\|^2. \end{aligned}$$

where δ_{B_k} denotes a change of variables restricted on $(\Delta\alpha_{B_k}, \Delta\mathbf{t}_{B_k})$ with all other variables fixed. Note the minimization in (69) is separable w.r.t $\{\delta_{B_k}\}_{k=1}^K$. Therefore, taking expectation of LHS and RHS w.r.t. k yields the result. \square

Before moving on, note that function $g(\mathbf{z}) = \frac{\eta_t}{2}\|\mathbf{z} + \mathbf{w}^t/\eta_t\|^2$ is locally Lipschitz-continuous with constant $L_g = \eta_t R_z$ for \mathbf{z} satisfying $\|\mathbf{z} + \mathbf{w}^t/\eta_t\| \leq R_z$, that is,

$$|g(\mathbf{z}_1) - g(\mathbf{z}_2)| \leq L_g \|\mathbf{z}_1 - \mathbf{z}_2\| \quad (37)$$

for $\forall \mathbf{z}_1, \mathbf{z}_2$ with $\|\mathbf{z}_1 + \mathbf{w}^t/\eta_t\| \leq R_z, \|\mathbf{z}_2 + \mathbf{w}^t/\eta_t\| \leq R_z$, where R_z is an upper bound on the magnitude of iterates $\|\mathbf{w}^{t+1}\|/\eta_t = \|\Phi^T \mathbf{x}^t + \mathbf{w}^t/\eta_t\|$.

From simplicity of analysis, in the following, we slightly loosen upper bounds by setting constants $L_g \leftarrow \max(L_g, 1)$, $M \leftarrow \max(M, 1)$, $\theta \leftarrow \max(\theta, 1)$, such that $L_g, M, \theta \geq 1$. Then we are ready to prove the main theorem of this section.

Theorem 5 (Linear Convergence). *The iterates $\{\mathbf{x}^s\}_{s=0}^\infty$ of Block Minimization for polyhedral $L_n(\cdot), R(\cdot)$ satisfy*

$$\mathbb{E}[F(\mathbf{x}^{s+1})] - F^* \leq \left(1 - \frac{1}{K\gamma}\right) (\mathbb{E}[F(\mathbf{x}^s)] - F^*)$$

where F^* is the optimum of (28) and

$$\gamma = \max\{16\eta_t M \theta (F^0 - F^*), 2M\theta(1 + 4L_g^2), 6\}.$$

Proof. Let $\mathbf{x}^* = \Pi_S(\mathbf{x}^s)$ be the projection of \mathbf{x}^s to the set of optimal solutions. From Lemma 3, we have

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}^{s+1})] - F(\mathbf{x}^s) &\leq \frac{1}{K} \left(\min_{\delta} h(\mathbf{x}^s + \delta) + \langle \nabla F(\mathbf{x}^s), \delta \rangle + \frac{M\eta_t}{2} \|\delta\|^2 \right) \\ &\leq \frac{1}{K} \left(\min_{\delta} h(\mathbf{x}^s + \delta) + F(\mathbf{x}^s + \delta) - F(\mathbf{x}^s) + \frac{M\eta_t}{2} \|\delta\|^2 \right) \\ &\leq \frac{1}{K} \left(\min_{a \in [0,1]} F(\mathbf{x}^s + a(\mathbf{x}^* - \mathbf{x}^s)) - F(\mathbf{x}^s) + \frac{M\eta_t a^2}{2} \|\mathbf{x}^* - \mathbf{x}^s\|^2 \right) \\ &\leq \frac{1}{K} \left(\min_{a \in [0,1]} -a(F(\mathbf{x}^s) - F(\mathbf{x}^*)) + \frac{M\eta_t a^2}{2} \|\mathbf{x}^* - \mathbf{x}^s\|^2 \right), \end{aligned} \quad (38)$$

where the second and fourth inequality follow from the convexity of $F(\mathbf{x})$, and the third inequality follows from the fact that both \mathbf{x}^* and \mathbf{x}^s are feasible ($h(\mathbf{x}^*) = h(\mathbf{x}^s) = 0$). Now based on the error bound inequality (68), we discuss two cases.

Case 1: $4L_g^2 \|\bar{\Phi}^T \mathbf{x} - \mathbf{p}^*\|^2 < (\mathbf{c}^T \mathbf{x} - q^*)^2$.

In this case, we have

$$\begin{aligned} \|\mathbf{x}^s - \mathbf{x}^*\|^2 &\leq \theta (\|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\|^2 + \|\mathbf{c}^T \mathbf{x}^s - q^*\|^2) \\ &\leq \theta \left(\frac{1}{4L_g^2} + 1 \right) (\mathbf{c}^T \mathbf{x}^s - q^*)^2 \leq 2\theta (\mathbf{c}^T \mathbf{x}^s - q^*)^2 \end{aligned} \quad (39)$$

and

$$|\mathbf{c}^T \mathbf{x}^s - q^*| \geq 2L_g \|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\| \geq 2|g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*)|.$$

Note in this case, $\mathbf{c}^T \mathbf{x}^s - q^*$ must be non-negative. Otherwise,

$$\begin{aligned} F(\mathbf{x}^s) - F^* &= g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*) + (\mathbf{c}^T \mathbf{x}^s - q^*) \\ &\leq |g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*)| - |\mathbf{c}^T \mathbf{x}^s - q^*| \\ &\leq -\frac{1}{2} |\mathbf{c}^T \mathbf{x}^s - q^*| < 0, \end{aligned}$$

leads to contradiction (since \mathbf{x}^s is feasible, $F(\mathbf{x}^s)$ cannot be smaller than F^*). Therefore, we have

$$\begin{aligned} F(\mathbf{x}^s) - F^* &= g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*) + \mathbf{c}^T \mathbf{x}^s - q^* \\ &\geq -|g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*)| + \mathbf{c}^T \mathbf{x}^s - q^* \\ &\geq \frac{1}{2}(\mathbf{c}^T \mathbf{x}^s - q^*). \end{aligned} \quad (40)$$

Combining (38), (39), and (40), we have

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}^{s+1})] - F(\mathbf{x}^s) &\leq \frac{1}{K} \min_{a \in [0,1]} -\frac{a}{2}(\mathbf{c}^T \mathbf{x}^s - q^*) + \frac{2\eta_t M \theta a^2}{2} (\mathbf{c}^T \mathbf{x}^s - q^*)^2 \\ &= \begin{cases} -1/(16\eta_t M \theta K) & , 1/(4\eta_t M \theta (\mathbf{c}^T \mathbf{x}^s - q^*)) \leq 1 \\ -\frac{1}{4K}(\mathbf{c}^T \mathbf{x}^s - q^*) & , o.w. \end{cases} \end{aligned}$$

Furthermore, we have

$$-\frac{1}{16\eta_t M \theta K} \leq -\frac{1}{16\eta_t M \theta K (F^0 - F^*)} (F(\mathbf{x}^*) - F^*)$$

where $F^0 = F(\mathbf{x}^0)$, and

$$-\frac{1}{4K}(\mathbf{c}^T \mathbf{x}^s - q^*) \leq -\frac{1}{6K}(F(\mathbf{x}^s) - F^*)$$

since $F(\mathbf{x}^s) - F^* \leq |g(\bar{\Phi}^T \mathbf{x}^s) - g(\mathbf{p}^*)| + \mathbf{c}^T \mathbf{x}^s - q^* \leq \frac{3}{2}(\mathbf{c}^T \mathbf{x}^s - q^*)$. In summary, for Case 1 we obtain

$$\mathbb{E}[F(\mathbf{x}^{s+1})] - F^* \leq (1 - \frac{1}{K\gamma_1}) (\mathbb{E}[F(\mathbf{x}^s)] - F^*) \quad (41)$$

where

$$\gamma_1 = \max \{ 16\eta_t M \theta (F^0 - F^*) , 6 \}. \quad (42)$$

Case 2: $4L_g^2 \|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\|^2 \geq (\mathbf{c}^T \mathbf{x}^s - q^*)^2$.

In this case, we have

$$\|\mathbf{x}^s - \mathbf{x}^*\|^2 \leq \theta (1 + 4L_g^2) \|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\|^2, \quad (43)$$

and by strong convexity of $g(\mathbf{z})$,

$$F(\mathbf{x}^s) - F^* \geq \mathbf{c}^T (\mathbf{x}^s - \mathbf{x}^*) + \nabla g(\mathbf{p}^*)^T \bar{\Phi}^T (\mathbf{x}^s - \mathbf{x}^*) + \frac{\eta_t}{2} \|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\|^2.$$

Adding inequality $0 = h(\mathbf{x}^s) - h(\mathbf{x}^*) \geq \langle \boldsymbol{\rho}^*, \mathbf{x}^s - \mathbf{x}^* \rangle$ for some $\boldsymbol{\rho}^* \in \partial h(\mathbf{x}^*)$ to the above gives

$$F(\mathbf{x}^s) - F^* \geq \frac{\eta_t}{2} \|\bar{\Phi}^T \mathbf{x}^s - \mathbf{p}^*\|^2 \quad (44)$$

since $\boldsymbol{\rho}^* + \mathbf{c} + \nabla g(\mathbf{p}^*)^T \bar{\Phi}^T = \boldsymbol{\rho}^* + \nabla F(\mathbf{x}^*) = 0$. Combining (38), (43), and (44), we obtain

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}^{s+1})] - F(\mathbf{x}^s) &\leq \frac{1}{K} \min_{a \in [0,1]} -a(F(\mathbf{x}^s) - F^*) + \frac{M\theta(1 + 4L_g^2)a^2}{2} (F(\mathbf{x}^s) - F^*) \\ &= -\frac{1}{2M\theta(1 + 4L_g^2)K} (F(\mathbf{x}^s) - F^*) \end{aligned} \quad (45)$$

Combining results of Case 1 (41) and Case 2 (45), and taking expectation on both sides w.r.t. the history leads to the result. \square

Theorem 2-2 (Dual-RBCD for Polyhedral Loss). *Let the sequence $\{\boldsymbol{\alpha}^s\}_{s=1}^\infty$ be the iterates produced by RBCD in the inner loop of Algorithm 2, and K be the number of blocks. Denote $\tilde{F}^*(\boldsymbol{\alpha})$ as the augmented dual objective function (18) and \tilde{F}_{opt}^* the optimum of (18). With probability $1 - \rho$,*

$$\tilde{F}^*(\boldsymbol{\alpha}^s) - \tilde{F}_{opt}^* \leq \epsilon, \text{ for } s \geq \gamma K \log\left(\frac{\tilde{F}^*(\boldsymbol{\alpha}^0) - \tilde{F}_{opt}^*}{\rho\epsilon}\right) \quad (46)$$

for some constant γ if $L_n(\cdot)$ and $R(\cdot)$ are polyhedral.

Proof. This simply applies Theorem 1 of [13] to transfer the linear convergence in expectation into high-probability iteration complexity. \square

6.2 Convergence of Proximal-Point Method

The proof of Theorem 3 comprises two parts. The first part proves linear convergence of Proximal-Point update under assumption that both loss $L_n(\cdot)$ and regularizer $R(\cdot)$ are either strictly convex and smooth or polyhedral. The second part proves a sublinear-type convergence depending on parameter η that holds for general convex function. The second part can be found in, for example, Theorem 2 of [18]. Here we prove the first part.

Here we prove linear convergence of ALM on problem (25) by leveraging some lemmas provided in the recent advance of analysis for Alternating Direction Method of Multiplier (ADMM) [5]. In particular, by taking Proximal-Point updates (or, equivalently, the ALM updates) as performing gradient descent on the convex, smooth function

$$G(\tilde{\mathbf{w}}) = \min_{\mathbf{w}} \sum_n L_n(\bar{\Phi}_n \mathbf{w}) + R(\mathbf{w}) + \frac{1}{2\eta} \|\mathbf{w} - \tilde{\mathbf{w}}\|^2 \quad (47)$$

and utilizing error bound proved in [5], we show that the Proximal-Point method linearly converges to the optimum of objective (25).

The following lemma establishes the fact $G(\tilde{\mathbf{w}})$ is smooth and its gradient $\nabla G(\tilde{\mathbf{w}})$ is Lipschitz continuous with modulus $\frac{1}{\eta}$.

Lemma 4. *The gradient of $G(\tilde{\mathbf{w}})$ is of the form*

$$\nabla G(\tilde{\mathbf{w}}) = -\left(\sum_{n=1}^N \Phi_n^T \alpha_n(\tilde{\mathbf{w}}) - \boldsymbol{\mu}(\tilde{\mathbf{w}})\right) \quad (48)$$

where $\alpha_n(\tilde{\mathbf{w}})$, $\boldsymbol{\mu}(\tilde{\mathbf{w}})$ are minimizers of (28). Furthermore, the gradient $\nabla G(\tilde{\mathbf{w}})$ is Lipschitz continuous with modulus $\frac{1}{\eta}$.

Proof. The convex objective function (25) fits the form of objective investigated in Multi-block ADMM [5]. Therefore, the theorem follows directly from Lemma 2.1, 2.2 of [5] respectively. \square

As a result of Lemma 4, the proximal-point update is exactly gradient descent of step size η , which when performed on a smooth function $G(\tilde{\mathbf{w}})$, guarantees descent amount

$$G(\mathbf{w}^{t+1}) - G(\mathbf{w}^t) \leq -\frac{\eta \|\nabla G(\mathbf{w}^t)\|^2}{2}. \quad (49)$$

The following theorem then guarantees linear convergence of ALM on our objective (25).

Theorem 6. *Denote S as the set of optimal solutions to (47) and $\Pi_S(\mathbf{w})$ as the projection of \mathbf{w} to S , and let G^* be the optimal function value. The iterates $\{\mathbf{w}^t\}_{t=1}^\infty$ produced by proximal-point method have*

$$\left\| \sum_{n=1}^N \Phi_n^T \alpha_n(\tilde{\mathbf{w}}) - \boldsymbol{\mu}(\tilde{\mathbf{w}}) \right\| = \|\nabla G(\mathbf{w}^t)\| \leq \epsilon$$

for number of iterations

$$t \geq \frac{4\tau}{\eta} \log\left(\sqrt{\frac{2(G(\mathbf{w}^0) - G^*)}{\eta}} \frac{1}{\epsilon}\right),$$

where $\tau > 0$ is a constant depending on S and initial distance to optimal set $\|\mathbf{w}^0 - \Pi_S(\mathbf{w}^0)\|$.

Proof. Since $L_n(\cdot)$ and $R(\cdot)$ are either strictly convex and smooth or polyhedral, $L_n^*(\cdot)$ and $R^*(\cdot)$ are also strictly convex and smooth or polyhedral. Therefore, problem (25) satisfies Assumption A(a)-A(e) of [5], and thus the error bound

$$G(\tilde{\mathbf{w}}) - G^* \leq \tau \|\nabla G(\tilde{\mathbf{w}})\|^2 \quad (50)$$

in Lemma 3.1 of [5] applies to $G(\tilde{\mathbf{w}})$ with compact domain $\tilde{\mathbf{w}} \in R(\mathbf{w}^0)$, where $\tau > 0$ is a constant that depends on geometry of S and the initial distance to the set of optimal solutions, and

$$R(\mathbf{w}^0) = \{\tilde{\mathbf{w}} \mid \|\tilde{\mathbf{w}} - \Pi_S(\tilde{\mathbf{w}})\| \leq \|\mathbf{w}^0 - \Pi_S(\mathbf{w}^0)\|\}.$$

is the set of $\tilde{\mathbf{w}}$ that lie within a radius of $\|\mathbf{w}^0 - \Pi_S(\mathbf{w}^0)\|$ to the set S . Note the iterates $\{\mathbf{w}^t\}_{t=0}^\infty$ all lie in the set $R(\mathbf{w}^0)$ by the non-expansiveness of proximal operation. Therefore, the error bound (50) applies to all iterates. Combining (69) and (50), we have

$$G(\mathbf{w}^{t+1}) - G(\mathbf{w}^t) \leq -\frac{\eta(G(\mathbf{w}^t) - G^*)}{2\tau},$$

which implies linear convergence. Let $\Delta G_t = G(\mathbf{w}^t) - G^*$, and we have

$$\Delta G_t \leq \left(1 - \frac{\eta}{2\tau}\right)^t \Delta G_0 \leq e^{-\frac{\eta t}{2\tau}} \Delta G_0 \leq \epsilon_1$$

when

$$t \geq \frac{2\tau}{\eta} \log\left(\frac{\Delta G_0}{\epsilon_1}\right).$$

Furthermore, by smoothness of $\nabla G(\cdot)$, we have

$$\Delta G_t \geq \frac{\eta \|\nabla G(\mathbf{w}^t)\|^2}{2}.$$

Therefore, to guarantee $\|\nabla G(\mathbf{w}^t)\| \leq \epsilon_2$, it suffices to have

$$\Delta G^t \leq \eta \epsilon_2^2 / 2,$$

which can be guaranteed by running

$$t \geq \frac{4\tau}{\eta} \log\left(\sqrt{\frac{2\Delta G_0}{\eta}} \frac{1}{\epsilon_2}\right)$$

iterations. □

Theorem 7 (Inexact Proximal Map). *Suppose, for a given dual iterate \mathbf{w}^t , each sub-problem (11) is solved inexactly s.t. the solution $\hat{\mathbf{w}}^{t+1}$ has*

$$\|\hat{\mathbf{w}}^{t+1} - \mathbf{prox}_{\eta_t F}(\mathbf{w}^t)\| \leq \epsilon_0. \quad (51)$$

Then let $\{\hat{\mathbf{w}}^t\}_{t=1}^\infty$ be the sequence of iterates produced by inexact proximal updates and $\{\mathbf{w}^t\}_{t=1}^\infty$ as that generated by exact updates. After t iterations, we have

$$\|\hat{\mathbf{w}}^t - \mathbf{w}^t\| \leq t\epsilon_0. \quad (52)$$

Proof. By the non-expansiveness of proximal operation,

$$\begin{aligned} \|\hat{\mathbf{w}}^{t+1} - \mathbf{w}^{t+1}\| &\leq \|\hat{\mathbf{w}}^{t+1} - \mathbf{prox}_{\eta_t F}(\hat{\mathbf{w}}^t)\| + \|\mathbf{prox}_{\eta_t F}(\hat{\mathbf{w}}^t) - \mathbf{w}^{t+1}\| \\ &\leq \epsilon_0 + \|\mathbf{prox}_{\eta_t F}(\hat{\mathbf{w}}^t) - \mathbf{prox}_{\eta_t F}(\mathbf{w}^t)\| \\ &\leq \epsilon_0 + \|\hat{\mathbf{w}}^t - \mathbf{w}^t\|. \end{aligned}$$

Recursively applying the above inequality leads to the conclusion. □

7 Appendix-B. ADMM under Limited Memory

In this section, we show how an algorithm for *distributed optimization* can be adapted for *limited-memory learning*, which then serves as a baseline to methods specially designed for limited-memory environment. In particular, the adaption sequentializes parallel computation performed on multiple machines into a series of tasks performed on single machine, where states and data partition of each simulated machine are loaded from (saved to) secondary storage units beforehand (afterward). As an example, we show how to adapt *Alternating Direction Method of Multiplier (ADMM)*, a recently proposed distributed optimization framework [1], into our setting.

Given a problem of the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^N f_i(\mathbf{w}), \quad (53)$$

Algorithm 3 ADMM (limited memory)

1. Split data \mathcal{D} into blocks B_1, B_2, \dots, B_K .
 2. Initialize $\mathbf{w}_k^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{0}, \boldsymbol{\mu}_k^0 = \mathbf{0}$.
 - for** $t = 0, 1, \dots$ (outer iteration) **do**
 3. $\mathbf{z}^{t+1} = \mathbf{0}$
 - for** $k = 1, 2, \dots, K$ **do**
 - 4.1. Swap data block $B_k, \mathbf{w}_k, \boldsymbol{\mu}_k$ into memory.
 - 4.2. $\mathbf{w}_k^{t+1} = \operatorname{argmin}_{\mathbf{w}} \mathcal{L}_k(\mathbf{w}, \mathbf{z}^t, \boldsymbol{\mu}_k^t)$
 - 4.3. $\mathbf{z}^{t+1} += (\mathbf{w}_k^{t+1} + \boldsymbol{\mu}_k^t / \rho) / K$
 - end for**
 - for** $k = 1, 2, \dots, K$ **do**
 - 5.1. Swap $\mathbf{w}_k^{t+1}, \boldsymbol{\mu}_k^t$ into memory.
 - 5.2. $\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t + \eta(\mathbf{w}_k^{t+1} - \mathbf{z}^{t+1})$.
 - end for**
 - end for**
-

Algorithm 4 Block-Coordinate ADMM (BC-ADMM)

1. Split data \mathcal{D} into blocks B_1, B_2, \dots, B_K .
 2. Initialize $\mathbf{w}_k^0 = \mathbf{0}, \mathbf{z}^0 = \mathbf{0}, \boldsymbol{\mu}_k^0 = \mathbf{0}$.
 - for** $t = 0, 1, \dots$ **do**
 - 3.1. Randomly chosen $k \in \{1..K\}$ w/o replacement.
 - 3.2. Swap data block $B_k, \mathbf{w}_k^t, \boldsymbol{\mu}_k^t$ into memory.
 - 3.3. $\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t + \eta(\mathbf{w}_k^t - \mathbf{z}^t)$.
 - 3.4. $\mathbf{w}_k^{t+1} = \operatorname{argmin}_{\mathbf{w}} \mathcal{L}_k(\mathbf{w}, \mathbf{z}^t, \boldsymbol{\mu}_k^{t+1})$
 - 3.5. $\mathbf{z}^{t+1} = \mathbf{z}^t + (\mathbf{w}_k^{t+1} + \boldsymbol{\mu}_k^{t+1} / \rho) / K - (\mathbf{w}_k^t + \boldsymbol{\mu}_k^t / \rho) / K$
 - end for**
-

the ADMM framework splits (53) into K smaller sub-problems defined on different data blocks B_1, B_2, \dots, B_K , and formulate the dual problem of

$$\begin{aligned} \min_{\mathbf{w}_k, \mathbf{z}} \quad & \sum_{k=1}^K f_k(\mathbf{w}_k) + \frac{\rho}{2} (\mathbf{w}_k - \mathbf{z})^2 \\ \text{s.t.} \quad & \mathbf{w}_k - \mathbf{z} = 0, \quad k = 1, \dots, K, \end{aligned} \quad (54)$$

where $f_k(\mathbf{w}_k) = \sum_{i \in B_k} f_i(\mathbf{w}_k)$, \mathbf{z} is the *consensus parameters*, and $\rho > 0$ is a hyper-parameter. The ADMM procedure finds the saddle point of Lagrangian

$$\max_{\boldsymbol{\mu}_k} \min_{\mathbf{w}_k, \mathbf{z}} \mathcal{L}(\mathbf{w}, \mathbf{z}, \boldsymbol{\mu}) = \sum_{k=1}^K f(\mathbf{w}_k) + \boldsymbol{\mu}_k^T (\mathbf{w}_k - \mathbf{z}) + \frac{\rho}{2} \|\mathbf{w}_k - \mathbf{z}\|^2 \quad (55)$$

via the following iterate

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \mathbf{z}^t, \boldsymbol{\mu}^t) \quad (56)$$

$$\mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z}} \mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}, \boldsymbol{\mu}^t) \quad (57)$$

$$\boldsymbol{\mu}_k^{t+1} = \boldsymbol{\mu}_k^t + \eta(\mathbf{w}_k^{t+1} - \mathbf{z}^{t+1}), \quad k = 1, \dots, K, \quad (58)$$

where η is a constant step size. Since given \mathbf{z}^t , $\mathcal{L}(\mathbf{w}, \mathbf{z}^t, \boldsymbol{\mu}^t)$ is separable w.r.t. $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K$, step (56) can be solved separately for each \mathbf{w}_k as

$$\mathbf{w}_k^{t+1} = \operatorname{argmin}_{\mathbf{w}_k} \mathcal{L}_k(\mathbf{w}_k, \mathbf{z}^t, \boldsymbol{\mu}_k^t), \quad k = 1, \dots, K. \quad (59)$$

Since the bottleneck of iterate lies in (59), ADMM is inherently suitable for distributed optimization via solving the K subproblems (59) on K machines. The only step requiring communication is (57),

which has close-form solution

$$\mathbf{z}^{t+1} = \frac{1}{K} \sum_{k=1}^K \mathbf{w}_k^{t+1} + \boldsymbol{\mu}_k^t / \rho, \quad (60)$$

that is, a simple average over parameters and multipliers. In limited-memory environment, however, only one block B_k of samples can be fit into memory at a time, and thus K times of swapping is required at each iteration. A naive implementation is depicted in Algorithm 3. Note, in some high-dimensional problem, the model parameters \mathbf{w}_k and $\boldsymbol{\mu}_k$ can be of comparable size to the data block, and thus need to be stored out of memory. One drawback of algorithm 3 is that the *consensus parameter* \mathbf{z} is not updated until K subproblems are solved. In Algorithm 4, we propose another adaption that updates the dual variables of one randomly chosen block B_k at a time as follows

$$\boldsymbol{\mu}_k^t = \boldsymbol{\mu}_k^{t-1} + \eta(\mathbf{w}_k^t - \mathbf{z}^t) \quad (61)$$

$$\mathbf{w}_k^{t+1} = \underset{\mathbf{w}_k}{\operatorname{argmin}} \mathcal{L}_k(\mathbf{w}_k, \mathbf{z}^t, \boldsymbol{\mu}_k^t) \quad (62)$$

$$\mathbf{z}^{t+1} = \underset{\mathbf{z}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w}_1^t, \dots, \mathbf{w}_k^{t+1}, \dots, \mathbf{w}_K^t, \mathbf{z}, \boldsymbol{\mu}_k^t). \quad (63)$$

In this version of limited-memory ADMM, the information learnt from one block can be passed to the next subproblem immediately, and consensus parameters $\boldsymbol{\mu}_k$, \mathbf{w}_k only need to be swapped once for each iteration. It has been shown that standard ADMM iterates in Algorithm 3 have global linear convergence to the optimum [5]. The following theorem shows the same type of convergence guarantee also applies to Algorithm 4 .

Theorem 8 (BC-ADMM Convergence). *Consider a regularized ERM problem (53) of the form*

$$f_i(\mathbf{w}) = L_i(\Phi_i \mathbf{w}) + \frac{1}{N} R(\mathbf{w}).$$

Let $d(\boldsymbol{\mu})$ be the dual function value of problem (54). If the loss function $L_i(\cdot)$ is smooth, and one of $L_i(\cdot)$ or $R(\cdot)$ is strongly convex, Algorithm 4 converges to the optimum of (53) at a linear rate, that is,

$$E[\Delta_p^t + \Delta_d^t] \leq \frac{1}{1 + \lambda} (\Delta_p^{t-1} + \Delta_d^{t-1}) \quad (64)$$

for some constant $\lambda > 0$, where

$$\begin{aligned} \Delta_p^t &= \mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - d(\boldsymbol{\mu}^t) \\ \Delta_d^t &= d^* - d(\boldsymbol{\mu}^t) \end{aligned} \quad (65)$$

are the primal and dual residuals at iterate t respectively.

Though being effective, the adapted algorithm takes little advantage of the sequential nature of limited-memory setting. In particular, since the distributed learning algorithm is designed to allow parallel updates, the information passed among parallel sub-problems is limited and the updates on dual variables (58), (62) are conservative with step size η compared to the exact block-coordinate minimization (12) in the Dual-Augmented Block Minimization framework. Note ADMM can be seen as an approximate Gradient Descent method on the dual, while analysis in coordinate descent literature [13] shows that Block-Coordinate descent can be up to K times faster than Gradient Descent in the worst-conditioned case, where K is the number of blocks.

8 Appendix-C. Convergence of Block-Coordinate ADMM

Let $d(\boldsymbol{\mu}) = \min_{\mathbf{w}, \mathbf{z}} \mathcal{L}(\mathbf{w}, \mathbf{z}, \boldsymbol{\mu})$ be the dual objective for $\boldsymbol{\mu}$ and $d^* = \max_{\boldsymbol{\mu}} d(\boldsymbol{\mu})$ be the optimal dual objective value, we define primal residual Δ_p^t and dual residual Δ_d^t of current iterate $(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t)$ as

$$\begin{aligned} \Delta_p^t &= \mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - d(\boldsymbol{\mu}^t) \\ \Delta_d^t &= d^* - d(\boldsymbol{\mu}^t). \end{aligned} \quad (66)$$

Note $\Delta_p^t \geq 0$, $\Delta_d^t \geq 0$, and $\Delta_p^t = \Delta_d^t = 0$ if and only if $(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t)$ is optimal.

Lemma 5 (Dual Iterate). For all $t \geq 1$,

$$\Delta_d^t - \Delta_d^{t-1} \leq -\eta(\mathbf{w}_k^t - \mathbf{z}^t)^T(\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t),$$

where $(\bar{\mathbf{w}}^t, \bar{\mathbf{z}}^t)$ is the solution to $\min_{\mathbf{w}, \mathbf{z}} \mathcal{L}(\mathbf{w}, \mathbf{z}, \boldsymbol{\mu}^t)$ that is closest to $(\mathbf{w}^t, \mathbf{z}^t)$.

Proof.

$$\begin{aligned} \Delta_d^t - \Delta_d^{t-1} &= d(\boldsymbol{\mu}^{t-1}) - d(\boldsymbol{\mu}^t) \\ &= \mathcal{L}(\bar{\mathbf{w}}^{t-1}, \bar{\mathbf{z}}^{t-1}, \boldsymbol{\mu}^{t-1}) - \mathcal{L}(\bar{\mathbf{w}}^t, \bar{\mathbf{z}}^t, \boldsymbol{\mu}^t) \\ &\leq \mathcal{L}(\bar{\mathbf{w}}^t, \bar{\mathbf{z}}^t, \boldsymbol{\mu}^{t-1}) - \mathcal{L}(\bar{\mathbf{w}}^t, \bar{\mathbf{z}}^t, \boldsymbol{\mu}^t) \\ &= (\boldsymbol{\mu}^{t-1} - \boldsymbol{\mu}^t)^T(\bar{\mathbf{w}}^t - \bar{\mathbf{z}}^t) \\ &= (\boldsymbol{\mu}_k^{t-1} - \boldsymbol{\mu}_k^t)^T(\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t) \\ &= -\eta(\mathbf{w}_k^t - \mathbf{z}^t)^T(\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t), \end{aligned}$$

where the third inequality follows from definition $(\bar{\mathbf{w}}^{t-1}, \bar{\mathbf{z}}^{t-1}) = \operatorname{argmin}_{\mathbf{w}, \mathbf{z}} \mathcal{L}(\mathbf{w}, \mathbf{z}, \boldsymbol{\mu}^{t-1})$. \square

Lemma 6 (Primal Iterate). For all $t \geq 1$,

$$\begin{aligned} \Delta_p^t - \Delta_p^{t-1} &\leq -\rho (\|\mathbf{w}_k^{t+1} - \mathbf{w}_k^t\|^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|^2) \\ &\quad + \eta (\|\mathbf{w}_k^t - \mathbf{z}^t\|^2 - (\mathbf{w}_k^t - \mathbf{z}^t)^T(\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t)) \end{aligned}$$

Proof.

$$\begin{aligned} \Delta_p^t - \Delta_p^{t-1} &= \\ &(\mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - d(\boldsymbol{\mu}^t)) - (\mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^{t-1}) - d(\boldsymbol{\mu}^{t-1})), \end{aligned}$$

where $d(\boldsymbol{\mu}^{t-1}) - d(\boldsymbol{\mu}^t)$ can be obtained via Lemma 2.1 as

$$d(\boldsymbol{\mu}^{t-1}) - d(\boldsymbol{\mu}^t) = -\eta(\mathbf{w}_k^t - \mathbf{z}^t)^T(\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t). \quad (67)$$

It remains to find

$$\begin{aligned} &\mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^{t-1}) = \\ &\mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t) + \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^{t-1}). \end{aligned}$$

From strong convexity of the augmented term $\frac{\rho}{2}\|\mathbf{w}_k - \mathbf{z}\|^2$, and that $\mathbf{w}^{t+1}, \mathbf{z}^{t+1}$ are minimizers for (56) and (57) respectively, we can bound the primal descent amount by

$$\begin{aligned} &\mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t) \\ &= \mathcal{L}_k(\mathbf{w}_k^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}_k^t) - \mathcal{L}_k(\mathbf{w}_k^t, \mathbf{z}^t, \boldsymbol{\mu}_k^t) \\ &\leq -\rho (\|\mathbf{w}_k^{t+1} - \mathbf{w}_k^t\|^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|^2). \end{aligned}$$

It is also known that

$$\mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^{t-1}) = \eta\|\mathbf{w}_k^t - \mathbf{z}^t\|^2.$$

Therefore,

$$\begin{aligned} &\mathcal{L}(\mathbf{w}^{t+1}, \mathbf{z}^{t+1}, \boldsymbol{\mu}^t) - \mathcal{L}(\mathbf{w}^t, \mathbf{z}^t, \boldsymbol{\mu}^{t-1}) \\ &\leq -\rho (\|\mathbf{w}_k^{t+1} - \mathbf{w}_k^t\|^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|^2) + \eta\|\mathbf{w}_k^t - \mathbf{z}^t\|^2. \end{aligned}$$

Combine above inequality with (67), we obtain the conclusion. \square

The following theorem guarantees descent of the primal-dual residual $\Delta_p^t + \Delta_d^t$ in expectation for each iteration of BC-ADMM.

Theorem 9 (Guaranteed Descent). For step-size η sufficiently small,

$$E[\Delta_p^t + \Delta_d^t] < (\Delta_p^{t-1} + \Delta_d^{t-1})$$

for all $t \geq 1$, where $E[\cdot]$ is expectation over blocks k_1, k_2, \dots, k_R drawn at iteration t .

Proof. Define

$$\Delta \mathbf{z}_k^t = (\mathbf{w}_k^{t+1} + \boldsymbol{\mu}_k^t / \rho) / K - (\mathbf{w}_k^t + \boldsymbol{\mu}_k^{t-1} / \rho) / K$$

and

$$\Delta \mathbf{z}^t = \frac{1}{K} \sum_{k=1}^K (\mathbf{w}_k^{t+1} + \boldsymbol{\mu}_k^t / \rho) - \frac{1}{K} \sum_{k=1}^K (\mathbf{w}_k^t + \boldsymbol{\mu}_k^{t-1} / \rho)$$

By Lemma 2.1 and 2.2, we have

$$\begin{aligned} & (\Delta_p^t + \Delta_d^t) - (\Delta_p^{t-1} + \Delta_d^{t-1}) \\ &= (\Delta_p^t - \Delta_p^{t-1}) + (\Delta_d^t - \Delta_d^{t-1}) \\ &\leq -\rho (\|\Delta \mathbf{w}_k^t\|^2 + \|\Delta \mathbf{z}_k^t\|^2) \\ &\quad + \eta (\|\mathbf{w}_k^t - \mathbf{z}^t\|^2 - 2(\mathbf{w}_k^t - \mathbf{z}^t)^T (\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t)). \end{aligned}$$

Taking expectation on both sides w.r.t. the random selected indexes k_1, k_2, \dots, k_R , we have

$$\begin{aligned} & E[\Delta_p^t + \Delta_d^t] - (\Delta_p^{t-1} + \Delta_d^{t-1}) \\ &\leq -\frac{\rho R}{K} (\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2) \\ &\quad + \frac{\eta R}{K} \left(\sum_{k=1}^K \|\mathbf{w}_k^t - \mathbf{z}^t\|^2 - 2 \sum_{k=1}^K (\mathbf{w}_k^t - \mathbf{z}^t)^T (\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t) \right), \end{aligned}$$

where $\Delta \mathbf{w}^t$ and $\Delta \mathbf{z}^t$ are the primal iterate of standard ADMM, and

$$\begin{aligned} & \sum_{k=1}^K \|\mathbf{w}_k^t - \mathbf{z}^t\|^2 - 2(\mathbf{w}_k^t - \mathbf{z}^t)^T (\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t) \\ &= \sum_{k=1}^K \|(\mathbf{w}_k^t - \mathbf{z}^t) - (\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t)\|^2 - \|\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t\|^2 \\ &\leq 2 \sum_{k=1}^K (\|\mathbf{w}_k^t - \bar{\mathbf{w}}_k^t\|^2 + \|\mathbf{z}^t - \bar{\mathbf{z}}^t\|^2) - \|\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t\|^2. \end{aligned}$$

Now we invoke the error bound in [5, Lemma 2.3, 2.5] to bound the distance between $(\mathbf{w}^t, \mathbf{z}^t)$ and $(\bar{\mathbf{w}}^t, \bar{\mathbf{z}}^t)$ in terms of progress in primal iterate $\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2$ as

$$\sum_{k=1}^K \|\mathbf{w}^t - \bar{\mathbf{w}}^t\|^2 + \|\mathbf{z}^t - \bar{\mathbf{z}}^t\|^2 \leq \tau (\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2), \quad (68)$$

where τ is a positive constant. Then we have

$$\begin{aligned} & E[\Delta_p^t + \Delta_d^t] - (\Delta_p^{t-1} + \Delta_d^{t-1}) \\ &\leq -\frac{R(\rho - 2\eta\tau)}{K} (\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2) - \frac{R\eta}{K} \sum_{k=1}^K \|\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t\|^2, \end{aligned} \quad (69)$$

which is always negative for $\eta < \rho / (2\tau)$. \square

Then we can have following theorem for linear convergence of BC-ADMM.

Theorem 10 (BC-ADMM Convergence). *Consider a regularized ERM problem (53) of the form*

$$f_i(\mathbf{w}) = L_i(\Phi_i \mathbf{w}) + \frac{1}{N} R(\mathbf{w}).$$

If the loss function $L_i(\cdot)$ is smooth, and one of $L_i(\cdot)$ or $R(\cdot)$ is strongly convex, Algorithm 4 converges to the optimum of (53) at a linear rate, that is,

$$E[\Delta_p^t + \Delta_d^t] \leq \frac{1}{1 + \lambda} (\Delta_p^{t-1} + \Delta_d^{t-1}) \quad (70)$$

for some constant $\lambda > 0$.

Proof. To prove linear convergence, we show that the two terms in (69) can be lower bounded by the current residual Δ_p^t, Δ_d^t respectively. In particular, we invoke the error bound in [5, Lemma 3.1] that shows

$$\Delta_d^t \leq \tau_2 \|\nabla d(\boldsymbol{\mu}^t)\| = \tau_2 \|\bar{\mathbf{w}}^t - \bar{\mathbf{z}}^t\|^2 \quad (71)$$

and

$$\Delta_p^t \leq \xi (\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2) \quad (72)$$

for some positive constant τ_2, ξ and $\forall t \geq t_0$, where (72) has combined [5, Lemma 3.1] and (68). Apply above error bounds on (69), we have

$$\begin{aligned} & E[\Delta_p^t + \Delta_d^t] - (\Delta_p^{t-1} + \Delta_d^{t-1}) \\ & \leq -\frac{R(\rho - 2\eta\tau)}{K} (\|\Delta \mathbf{w}^t\|^2 + \|\Delta \mathbf{z}^t\|^2) - \frac{R\eta}{K} \sum_{k=1}^K \|\bar{\mathbf{w}}_k^t - \bar{\mathbf{z}}^t\|^2 \\ & \leq -\frac{R(\rho - 2\eta\tau)}{K\xi} \Delta_p^t - \frac{R\eta}{K\tau_2} \Delta_d^t \\ & \leq -\lambda(\Delta_p^t + \Delta_d^t), \end{aligned}$$

where $\lambda = \frac{R}{K} \min\{(\rho - 2\eta\tau)\xi^{-1}, \eta\tau_2^{-1}\} > 0$ for step size $\eta < \rho/(2\tau)$. After rearrangement we have

$$E[\Delta_p^t + \Delta_d^t] \leq \frac{1}{1 + \lambda} (\Delta_p^{t-1} + \Delta_d^{t-1}), \quad t \geq t_0.$$

□