

A Omitted Proofs

A.1 Derivation of Primal-Dual Formulation

$$\begin{aligned}
\min_{\mathbf{x} \in C} P(\mathbf{x}) &= \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x}) \\
&= \min_{\mathbf{x} \in C, \mathbf{b} = A\mathbf{x}} \frac{1}{n} \sum_{i=1}^n f_i(b_i) + g(\mathbf{x}) \\
&= \min_{\mathbf{x} \in C, \mathbf{b}} \max_{\mathbf{y}} \left\{ \frac{1}{n} \sum_{i=1}^n f_i(b_i) + g(\mathbf{x}) + \frac{1}{n} \langle \mathbf{y}, A\mathbf{x} - \mathbf{b} \rangle \right\} \\
&= \min_{\mathbf{x} \in C} \max_{\mathbf{y}} \left\{ g(\mathbf{x}) + \frac{1}{n} \langle \mathbf{y}, A\mathbf{x} \rangle + \min_{\mathbf{b}} \left\{ \frac{1}{n} \sum_{i=1}^n f_i(b_i) - \frac{1}{n} \langle \mathbf{y}, \mathbf{b} \rangle \right\} \right\} \\
&= \min_{\mathbf{x} \in C} \max_{\mathbf{y}} \left\{ \mathcal{L}(\mathbf{x}, \mathbf{y}) := g(\mathbf{x}) + \frac{1}{n} \langle \mathbf{y}, A\mathbf{x} \rangle - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) \right\} \\
&= \max_{\mathbf{y}} \left\{ D(\mathbf{y}) := \min_{\mathbf{x} \in C} \left\{ g(\mathbf{x}) + \frac{1}{n} \langle \mathbf{y}, A\mathbf{x} \rangle - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) \right\} \right\}
\end{aligned}$$

We use Von Neumann-Fan minimax theorem for the whole derivation when swapping each min-max formula [7]. For the last equality, there is a convex constraint in the minimization part. Although the original Von Neumann-Fan doesn't have constraints, it naturally applies to the case when \mathbf{x} (assuming function is convex to \mathbf{x}) is bounded in a convex set, since we could change $f(\mathbf{x}, \mathbf{y})$ to $f(\mathbf{x}, \mathbf{y}) + I_C(\mathbf{x})$, where $I_C(\mathbf{x}) = 0$ if $\mathbf{x} \in C$ and ∞ otherwise. Then the property will be properly inherited.

A.2 Notation and simple facts

Recall primal, dual and Lagrangian forms:

$$\begin{aligned}
P(\mathbf{x}) &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{a}_i^\top \mathbf{x}) + g(\mathbf{x}) \\
\mathcal{L}(\mathbf{x}, \mathbf{y}) &\stackrel{\text{def}}{=} g(\mathbf{x}) + \frac{1}{n} \mathbf{y}^\top A\mathbf{x} - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) \\
D(\mathbf{y}) &\stackrel{\text{def}}{=} \min_{\mathbf{x} \in C} \mathcal{L}(\mathbf{x}, \mathbf{y}) = \mathcal{L}(\bar{\mathbf{x}}(\mathbf{y}), \mathbf{y})
\end{aligned}$$

Similar to the definitions in [24], we introduce the primal gap defined as $\Delta_p^{(t)} \stackrel{\text{def}}{=} \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - D(\mathbf{y}^{(t)})$, and dual gap $\Delta_d^{(t)} \stackrel{\text{def}}{=} D^* - D(\mathbf{y}^{(t)})$. Recall the assumptions:

- f_i is convex and β -smooth, and is α strongly convex over some convex set, and linear otherwise.
- $R = \max_i \|\mathbf{a}_i\|_2^2, \forall i \in [n]$.
- g is μ -strongly convex and L -smooth.

To begin with, it is easy to verify that f_i^* is $1/\beta$ -strongly convex and is $1/\alpha$ -smooth on a convex set and infinity otherwise (See Claim A.6). For simplicity we first assume $\alpha \geq \frac{1}{2}\beta$ and then generalize the result.

Claim A.1. • Since $D(\mathbf{y}) = \min_{\mathbf{x} \in C} \{g(\mathbf{x}) + \frac{1}{n} \mathbf{y}^\top A\mathbf{x}\} - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i)$, $-D(\mathbf{y})$ is at least $\frac{1}{\beta}$ -strongly convex.

- Based on our update rule, $\exists \mathbf{g} \in \partial_{\mathbf{y}} \frac{1}{n} \sum_i f_i^*(\mathbf{y}^{(t)})$, such that

$$\mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} = \delta \left(\frac{1}{n} A_{I^{(t)},:} \mathbf{x}^{(t)} - \mathbf{g}_{I^{(t)}} \right). \quad (19)$$

And our update rule ensures that $I^{(t)}$ consists of indices $i \in [n]$ that maximizes $|\frac{1}{n}\mathbf{a}_i^\top \mathbf{x}^{(t)} - g_i|$.

A.3 Primal Progress

Lemma A.2. (Primal Progress)

$$\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) \leq (1 - \frac{\eta}{2}) \left(\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) \right)$$

Or equivalently,

$$(1 - \frac{\eta}{2})(\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) \leq -\frac{\eta}{2} \left(\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) \right) \equiv -\frac{\eta}{2} \Delta_p^{(t)}$$

Proof. Simply replace h_t as $\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - D(\mathbf{y}^{(t)})$ and h_{t+1} as $\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - D(\mathbf{y}^{(t)})$ in Inequality (4). We could conclude that $h_{t+1} \leq (1 - \eta + \eta^2 \frac{L}{\mu})h_t$. Therefore when $\eta \leq \frac{\mu}{2L}$, $h_{t+1} \leq (1 - \frac{\eta}{2})h_t$ and the first part of Lemma A.2 is true. Some simple rearrangement suffices the second part of the lemma. \square

A.4 Primal Dual Progress

In order to get a clue on how to analyze the dual progress, we first look at how the primal and dual evolve through iterations.

For an index set I and a vector $\mathbf{y} \in \mathbb{R}^n$, denote $\mathbf{y}_I = \sum_{i \in I} y_i \mathbf{e}_i \in \mathbb{R}^k$ as the subarray of \mathbf{y} indexed by I , with $|I| = k$. Recall Algorithm 1 selects the coordinates to update in the dual variable as $I^{(t)}$.

Lemma A.3. (Primal-Dual Progress).

$$\begin{aligned} & \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\ & \leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{1}{2\delta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\ & \quad + \frac{2\delta Rk}{n^2} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|^2. \end{aligned}$$

Proof. Notice we have claimed that $-D(\mathbf{y})$ is $\frac{1}{\beta}$ -strongly convex and for all $\mathbf{g} \in \partial_{\mathbf{y}} \frac{1}{n} \sum_i f_i^*(\mathbf{y}^{(t)})$,

$$\begin{aligned} & \Delta_d^{(t)} - \Delta_d^{(t-1)} = (-D(\mathbf{y}^{(t)})) - (-D(\mathbf{y}^{(t-1)})) \\ & \leq \langle -\nabla_{\mathbf{y}} \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}), \mathbf{y}^{(t)} - \mathbf{y}^{(t-1)} \rangle - \frac{1}{2\beta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\ & = -\langle \frac{1}{n} A_{I^{(t)},:} \bar{\mathbf{x}}^{(t)} - \mathbf{g}_{I^{(t)}}, \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \rangle - \frac{1}{2\beta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \end{aligned} \quad (20)$$

Meanwhile since $-\mathcal{L}(\mathbf{x}, \mathbf{y})$ is $\frac{1}{\alpha}$ -smooth over its feasible set,

$$\begin{aligned} & \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t-1)}) \\ & = -\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t-1)}) - (-\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) \\ & \leq (\frac{1}{n} A_{I^{(t)},:} \mathbf{x}^{(t)} - \mathbf{g}_{I^{(t)}})^\top (\mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)}) + \frac{1}{2\alpha} \|\mathbf{y}_{I^{(t)}}^{(t-1)} - \mathbf{y}_{I^{(t)}}^{(t)}\|^2 \\ & = (\frac{1}{\delta} + \frac{1}{2\alpha}) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2. \end{aligned} \quad (21)$$

Also, with the update rule of dual variables, we could make use of Eqn. (19) and re-write Eqn. (20) as:

$$\begin{aligned} & \Delta_d^{(t)} - \Delta_d^{(t-1)} \\ & \leq -\langle \frac{1}{n} A_{I^{(t)},:} \bar{\mathbf{x}}^{(t)} - \mathbf{g}_{I^{(t)}}, \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \rangle - \frac{1}{\delta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\ & \quad + (\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)})^\top (\frac{1}{n} A_{I^{(t)},:} \mathbf{x}^{(t)} - \mathbf{g}_{I^{(t)}}) - \frac{1}{2\beta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\ & = -\langle \frac{1}{n} A_{I^{(t)},:} (\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}), \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \rangle - (\frac{1}{\delta} + \frac{1}{2\beta}) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \end{aligned} \quad (22)$$

Together we get:

$$\begin{aligned}
& \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
&= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t-1)}) \\
&\quad + 2(\Delta_d^{(t)} - \Delta_d^{(t-1)}) \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \left(\frac{1}{\delta} + \frac{1}{2\alpha}\right) \|\mathbf{y}_{I^{(t)}}^{(t-1)} - \mathbf{y}_{I^{(t)}}^{(t)}\|^2 + 2(\Delta_d^{(t)} - \Delta_d^{(t-1)}) \\
&\quad \text{(from Eqn. (21))} \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \left(\frac{1}{\delta} + \frac{1}{2\alpha}\right) \|\mathbf{y}_{I^{(t)}}^{(t-1)} - \mathbf{y}_{I^{(t)}}^{(t)}\|^2 \\
&\quad - 2\left\langle \frac{1}{n} A_{I^{(t)}}, (\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}), \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \right\rangle - 2\left(\frac{1}{\delta} + \frac{1}{2\beta}\right) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\quad \text{(from Eqn. (22))} \\
&= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - 2\left\langle \frac{1}{n} A_{I^{(t)}}, (\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}), \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \right\rangle \\
&\quad - \left(\frac{1}{\delta} + \frac{1}{\beta} - \frac{1}{2\alpha}\right) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + 2\delta \left\| \frac{1}{n} A_{I^{(t)}}, (\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}) \right\|^2 \\
&\quad - \left(\frac{1}{\delta} - \frac{1}{2\delta}\right) \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \quad \text{(since } 2ab \leq \gamma a^2 + 1/\gamma b^2) \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{1}{2\delta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\quad + \frac{2\delta Rk}{n^2} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|^2
\end{aligned}$$

□

Therefore we will connect the progress induced by $-\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|$ and dual gap $\Delta_d^{(t)}$ next.

A.5 Dual progress

Claim A.4. An α -strongly convex function f satisfies:

$$f(\mathbf{x}) - f^* \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|_2^2$$

This simply due to $f(\mathbf{x}) - f^* \leq \langle \nabla f(\mathbf{x}), \mathbf{x} - \bar{\mathbf{x}} \rangle - \frac{\alpha}{2} \|\mathbf{x} - \bar{\mathbf{x}}\|_2^2 \leq \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|^2 + \frac{\alpha}{2} \|\mathbf{x} - \bar{\mathbf{x}}\|^2 - \frac{\alpha}{2} \|\mathbf{x} - \bar{\mathbf{x}}\|^2 = \frac{1}{2\alpha} \|\nabla f(\mathbf{x})\|^2$.

Since $-D$ is $\frac{1}{\beta}$ -strongly convex, we get

$$\begin{aligned}
\Delta_d^{(t)} &= D^* - D(\mathbf{y}^{(t)}) \leq \frac{\beta}{2} \|\nabla D(\mathbf{y}^{(t)})\|_2^2 \\
&= \frac{\beta}{2} \left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \mathbf{g} \right\|_2^2 \\
&\leq \frac{n\beta}{2k} \left\| \frac{1}{n} A_{\bar{I}}, \bar{\mathbf{x}}^{(t)} - \mathbf{g}_{\bar{I}} \right\|_2^2, \tag{23}
\end{aligned}$$

where \bar{I} is a set of size k that maximizes the values of $A_{\bar{I}}^T \bar{\mathbf{x}}^{(t)} - \mathbf{g}_{\bar{I}}$.

Lemma A.5 (Dual Progress).

$$-\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \leq -\frac{k\delta}{n\beta} \Delta_d^{(t)} + \frac{k\delta}{n^2} R \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_2^2$$

Proof of Lemma A.5. Define $\Delta = \frac{1}{n}A(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)})$. Since

$$\begin{aligned}
& - \left\| \frac{1}{n}A_{I^{(t)}}^\top \mathbf{x}^{(t)} - \mathbf{g}_{I^{(t)}} \right\|^2 \\
& \leq - \left\| \frac{1}{n}A_I^\top \mathbf{x}^{(t)} - \mathbf{g}_I \right\|^2 && (\text{choice of } I^{(t)}) \\
& = - \left\| \frac{1}{n}A_I^\top \bar{\mathbf{x}}^{(t)} - \mathbf{g}_I - \Delta_I \right\|^2 \\
& \leq - \frac{1}{2} \left\| \frac{1}{n}A_I^\top \bar{\mathbf{x}}^{(t)} - \mathbf{g}_I \right\|^2 + \|\Delta_I\|_2^2 \\
& && (\text{since } -(a+b)^2 \leq -1/2a^2 + b^2) \\
& \leq - \frac{k}{n\beta} \Delta_d^{(t)} + \|\Delta_I\|_2^2 && (\text{from (23)}) \\
& \leq - \frac{k}{n\beta} \Delta_d^{(t)} + \frac{k}{n^2} R \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_2^2
\end{aligned}$$

With the relation between $\frac{1}{n}A_{I^{(t)}}^\top \mathbf{x}^{(t)} - \mathbf{g}_{I^{(t)}}$ and $\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}$ we finish the proof. \square

A.6 Convergence on Duality Gap

Now we are able to merge the primal/dual progress to get the overall progress on the duality gap.

Proof of Theorem 4.1. We simply blend Lemma A.2 and Lemma A.5 with the primal-dual progress (Lemma A.3):

$$\begin{aligned}
& \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
& \leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{1}{2\delta} \|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
& \quad + \frac{2\delta Rk}{n^2} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|^2 && (\text{Lemma A.3}) \\
& \leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \frac{\delta}{2} \left(-\frac{k}{n\beta} \Delta_d^{(t)} + \frac{k}{n^2} R \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_2^2 \right) \\
& \quad + \frac{2\delta Rk}{n^2} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|^2 && (\text{Lemma A.5}) \\
& = \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{k\delta}{2n\beta} \Delta_d^{(t)} + \frac{5R\delta k}{2n^2} \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_2^2 \\
& \leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{k\delta}{2n\beta} \Delta_d^{(t)} + \frac{5R\delta k}{\mu n^2} (\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})) \\
& = \left(1 - \frac{5R\delta k}{\mu n^2}\right) (\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) - \frac{k\delta}{2n\beta} \Delta_d^{(t)} \\
& \quad + \frac{5R\delta k}{\mu n^2} (\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})) \\
& \leq - \frac{k\delta}{2n\beta} \Delta_d^{(t)} - \left(\left(1 - \frac{5R\delta k}{\mu n^2}\right) \frac{\mu}{4L} - \frac{5R\delta k}{\mu n^2} \right) \Delta_p^{(t)} && (\text{Lemma A.2})
\end{aligned}$$

When setting $\frac{k\delta}{2n\beta} = \left(1 - \frac{5R\delta k}{\mu n^2}\right) \frac{\mu}{4L} - \frac{5R\delta k}{\mu n^2}$, we get that $\Delta^{(t)} \leq \frac{1}{1+a} \Delta^{(t-1)}$, where $1/a = \mathcal{O}(\frac{L}{\mu}(1 + \frac{R\beta}{n\mu}))$. Therefore it takes $\mathcal{O}(\frac{L}{\mu}(1 + \frac{R\beta}{n\mu}) \log \frac{1}{\epsilon})$ for $\Delta^{(t)}$ to reach ϵ .

When $\beta > 2\alpha$, we could redefine the primal-dual process as $\Delta^{(t)} := (\frac{\beta}{\alpha} - 1) \Delta_d^{(t)} + \Delta_p^{(t)}$ and rewrite some of the key steps, especially for the overall primal-dual progress.

$$\begin{aligned}
& \Delta^{(t)} - \Delta^{(t-1)} \\
&= \left(\frac{\beta}{\alpha} - 1\right)(\Delta_d^{(t)} - \Delta_d^{(t-1)}) + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
&= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t-1)}) \\
&\quad + \frac{\beta}{\alpha}(\Delta_d^{(t)} - \Delta_d^{(t-1)}) \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \left(\frac{1}{\delta} + \frac{1}{2\alpha}\right)\|\mathbf{y}_{I^{(t)}}^{(t-1)} - \mathbf{y}_{I^{(t)}}^{(t)}\|^2 \\
&\quad - \frac{\beta}{\alpha}\left\langle \frac{1}{n}A_{I^{(t)},:}(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}), \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \right\rangle - \frac{\beta}{\alpha}\left(\frac{1}{\delta} + \frac{1}{2\beta}\right)\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\hspace{15em} \text{(from Eqn. (21) and (22))} \\
&= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{\beta}{\alpha}\left\langle \frac{1}{n}A_{I^{(t)},:}(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}), \mathbf{y}_{I^{(t)}}^{(t)} - \mathbf{y}_{I^{(t)}}^{(t-1)} \right\rangle \\
&\quad - \left(\frac{\beta}{\alpha} - 1\right)\frac{1}{\delta}\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \frac{3\beta}{2\alpha}\delta\left\|\frac{1}{n}A_{I^{(t)},:}(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)})\right\|^2 \\
&\quad - \left(\frac{3\beta}{4\alpha} - 1\right)\frac{1}{\delta}\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \hspace{10em} \text{(since } ab \leq \delta a^2 + 1/(4\delta)b^2) \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) + \frac{\beta}{\alpha}\delta\left\|\frac{1}{n}A_{I^{(t)},:}(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)})\right\|^2 \\
&\quad - \frac{\beta}{4\alpha\delta}\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \hspace{10em} \text{(since } \beta/\alpha \geq 2) \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{\beta}{4\alpha\delta}\|\mathbf{y}^{(t)} - \mathbf{y}^{(t-1)}\|^2 \\
&\quad + \frac{\beta\delta Rk}{\alpha n^2}\|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|^2
\end{aligned}$$

Similarly to the previous setting, we get the whole primal-dual progress is bounded as follows:

$$\begin{aligned}
& \left(\frac{\beta}{\alpha} - 1\right)(\Delta_d^{(t)} - \Delta_d^{(t-1)}) + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
&\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \frac{\beta\delta}{4\alpha}\frac{k}{n\beta}\Delta_d^{(t)} \\
&\quad + \frac{5\beta R\delta k}{2\alpha\mu n^2}(\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})) \\
&\leq -\frac{\beta}{4\alpha}\frac{k\delta}{n\beta}\Delta_d^{(t)} - \left(1 - \frac{5\beta R\delta k}{2\alpha\mu n^2}\right)\frac{\mu}{4L} - \frac{5\beta R\delta k}{2\alpha\mu n^2}\Delta_p^{(t)}
\end{aligned}$$

Therefore, when we set a proper k and δ such that $\frac{\beta}{4\alpha}\frac{k\delta}{n\beta} = \left(\frac{\beta}{\alpha} - 1\right)\left(1 - \frac{5\beta R\delta k}{2\alpha\mu n^2}\right)\frac{\mu}{4L} - \frac{5\beta R\delta k}{2\alpha\mu n^2}$, and since $\frac{\beta}{\alpha} - 1 \geq \frac{\beta}{2\alpha}$, we get $\delta = \frac{1}{k}\left(\frac{L}{\mu n\beta} + \frac{5\beta R}{2\alpha\mu n^2}(1 + 4\frac{L}{\mu})\right)^{-1}$. And we have $\Delta^{(t)} - \Delta^{(t-1)} \leq -1/a\Delta^{(t)}$, where $a = \mathcal{O}(\frac{L}{\mu}(1 + \frac{\beta}{\alpha}\frac{R\beta}{n\mu}))$. Therefore it takes $t = \mathcal{O}(\frac{L}{\mu}(1 + \frac{\beta}{\alpha}\frac{R\beta}{n\mu})\log \frac{1}{\epsilon})$ iterations for the duality gap $\Delta^{(t)}$ to reach ϵ error. \square

A.7 Smooth Hinge Loss and Relevant Properties

Smooth hinge loss is defined as follows:

$$h(z) = \begin{cases} \frac{1}{2} - z & \text{if } z < 0 \\ \frac{1}{2}(1 - z)^2 & \text{if } z \in [0, 1] \\ 0 & \text{otherwise.} \end{cases} \quad (24)$$

Our loss function over a prediction p associated with a label $\ell_i \in \{\pm 1\}$ will be $f_i(p) = h(p\ell_i)$. The derivative of smooth hinge loss h is:

$$h'(z) = \begin{cases} -1 & \text{if } z < 0 \\ z - 1 & \text{if } z \in [0, 1] \\ 0 & \text{otherwise.} \end{cases} \quad (25)$$

Its convex conjugate is:

$$h^*(z^*) = \begin{cases} \frac{1}{2}(z^*)^2 + z^* & \text{if } z^* \in [-1, 0] \\ \infty & \text{otherwise.} \end{cases} \quad (26)$$

Notice since $f_i(p) = h(p\ell_i)$, $f_i^*(p) = h^*(p/\ell_i) = h^*(p\ell_i)$.

Claim A.6. *For a convex and β -smooth scalar function f , if it is α strongly convex over some convex set, and linear otherwise, then its conjugate function f^* is $1/\beta$ -strongly convex, and it is a $1/\alpha$ -smooth function plus an indicator function over some interval $[a, b]$.*

Proof. To begin with, since $f''(x) \leq \beta, \forall x$, meaning f is β -smooth, then with duality we have f^* is $1/\beta$ strongly convex [16]. Secondly, since f is α strongly convex over a convex set, meaning an interval for \mathbb{R} , therefore f could only be linear on $(-\infty, a]$ or $[b, \infty)$, and is α -strongly convex over the set $[a, b]$ (Here for simplicity $a < b$ could be $\pm\infty$). We denote $f'(-\infty) := \lim_{x \rightarrow -\infty} f'(x)$ and $f'(\infty)$ likewise. It's easy to notice that $f'(-\infty) \leq f'(a) < f'(b) \leq f'(\infty)$ since f is convex overall and strongly convex over $[a, b]$. Therefore $f(y) > f(a) + f'(a)(y - a)$ when $y > a$ and $f(y) = f(a) + f'(a)(y - a)$ when $y \leq a$.

Now since $f^*(x^*) \equiv \max_x \{x^*x - f(x)\}$, it's easy to verify that when $x^* < f'(a)$, $x^*x - f(x) = x^*x - f(a) - f'(a)(x - a) = -(f'(a) - x^*)x - f(a) + f'(a)a \rightarrow \infty$ when $x \rightarrow -\infty$. Similarly, when $x^* > f'(b)$, $f^*(x^*) = \infty$. On the other hand, when $x^* \in [f'(a), f'(b)]$, $f^*(x^*) = \max_x \{x^*x - f(x)\} = \max_{x \in [a, b]} \{x^*x - f(x)\}$. This is because $x^*a - f(a) \geq x^*y - f(y) = x^*y - f(y) - f'(a)(y - a), \forall y \leq a$, and similarly $x^*b - f(b) \geq x^*y - f(y), \forall y > b$. Therefore f^* is $1/\alpha$ smooth over the interval $[f'(a), f'(b)]$, where $-\infty \leq f'(a) < f'(b) \leq \infty$. □

A.8 Convergence of Optimization over Trace Norm Ball

The convergence analysis for trace norm ball is mostly similar to the case of ℓ_1 ball. The most difference lies on the primal part, where our approximated update incur linear progress as well as some error.

Lemma A.7 (Primal Progress for Algorithm 2). *Suppose $\text{rank } \bar{X}^{(t)} \leq s$ and $\epsilon > 0$. If each \tilde{X} computed in our algorithm is a $(\frac{1}{2}, \frac{\epsilon}{8})$ -approximate solution to (18), then for every t , it satisfies $\mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(X^{(t)}, Y^{(t)}) \leq -\frac{\mu}{8L} \Delta_p^{(t)} + \frac{\epsilon}{16}$.*

Proof. Refer to the proof in [1] we have:

$$\mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(\bar{X}^{(t)}, Y^{(t)}) \leq (1 - \frac{\mu}{8L}) (\mathcal{L}(X^{(t)}, Y^{(t)}) - \mathcal{L}(\bar{X}^{(t)}, Y^{(t)})) + \frac{\epsilon\mu}{16L}$$

Now move the first term on the RHS to the left and rearrange we get:

$$(1 - \frac{\mu}{8L})(\mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(X^{(t)}, Y^{(t)})) + \frac{\mu}{8L} (\mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(\bar{X}^{(t)}, Y^{(t)})) \leq \frac{\epsilon\mu}{16L}$$

Therefore we get:

$$\mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(X^{(t)}, Y^{(t)}) \leq -\frac{\mu}{8L} \Delta_p^{(t)} + \frac{\epsilon}{16}. \quad \square$$

Now back to the convergence guarantees on the trace norm ball.

Proof of Theorem 4.3. We again define $\Delta = \frac{1}{n}A(\bar{X}^{(t)} - X^{(t)})$. $G = \nabla_Y \mathcal{L}(X^{(t)}, Y^{(t)})$ such that $Y_{I^{(t)},:}^{(t)} - Y_{I^{(t)},:}^{(t-1)} = \delta(\frac{1}{n}\langle A_{I^{(t)},:}, X^{(t)} \rangle - G_{I^{(t)},:})$. Again we get $\|\Delta\|_F^2 \leq \frac{R}{n^2} \|\bar{X}^{(t)} - X^{(t)}\|_F^2$.

$$\Delta_d^{(t)} \leq \frac{\beta}{2} \left\| \frac{1}{n} A \bar{X}^{(t)} - G \right\|_F^2 \leq \frac{n\beta}{2k} \left\| \frac{1}{n} A_{I^{(t)},:} \bar{X}^{(t)} - G_{I^{(t)},:} \right\|_F^2$$

Other parts are exactly the same and we get:

$$\begin{aligned}
& \left(\frac{\beta}{\alpha} - 1\right)(\Delta_d^{(t)} - \Delta_d^{(t-1)}) + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
& \leq \mathcal{L}(X^{(t+1)}, Y^{(t)}) - \mathcal{L}(X^{(t)}, Y^{(t)}) - \frac{\beta\delta}{4\alpha} \frac{k}{n\beta} \Delta_d^{(t)} \\
& \quad + \frac{5\beta R\delta k}{2\alpha\mu n^2} (\mathcal{L}(X^{(t)}, Y^{(t)}) - \mathcal{L}(\bar{X}^{(t)}, Y^{(t)})) \\
& \leq -\frac{\beta}{4\alpha} \frac{k\delta}{n\beta} \Delta_d^{(t)} - \left(\left(1 - \frac{5\beta R\delta k}{2\alpha\mu n^2}\right) \frac{\mu}{8L} - \frac{5\beta R\delta k}{2\alpha\mu n^2} \right) \Delta_p^{(t)} + \left(1 - \frac{5\beta R\delta k}{2\alpha\mu n^2}\right) \frac{\epsilon}{16}
\end{aligned}$$

(Lemma A.7)

Therefore when $\delta \leq \frac{1}{k} \left(\frac{L}{\mu n\beta} + \frac{5\beta R}{2\alpha\mu n^2} (1 + 8\frac{L}{\mu}) \right)^{-1}$, it satisfies $\Delta^{(t)} - \Delta^{(t-1)} \leq -\frac{k\delta}{2\beta n} \Delta^{(t)} + \frac{\epsilon}{16}$. Therefore denote $a = \frac{2\beta n}{k\delta}$, we get $\Delta^{(t)} \leq \frac{a}{a+1} (\Delta^{(t-1)} + \frac{\epsilon}{16})$. Therefore we get $\Delta^{(t)} \leq \left(\frac{a}{a+1}\right)^t \Delta^{(0)} + \frac{\epsilon}{16} \sum_{i=1}^t \left(\frac{a}{a+1}\right)^i \leq \left(\frac{c}{c+1}\right)^t \Delta^{(0)} + \epsilon/16$. Since $\left(\frac{a}{a+1}\right)^t \leq e^{-t/a}$, it takes around $a = \mathcal{O}\left(\frac{L}{\mu} \left(1 + \frac{\beta R\beta}{\alpha n\mu}\right) \log \frac{1}{\epsilon}\right)$ iterations for the duality gap to get ϵ -error. \square

A.9 Difficulty on Extension to Polytope Constraints

Another important type of constraint we have not explored in this paper is the polytope constraint. Specifically,

$$\min_{\mathbf{x} \in M \subset \mathbb{R}^d} f(A\mathbf{x}) + g(\mathbf{x}), M = \text{conv}(\mathcal{A}), \text{ with only access to: } \text{LMO}_{\mathcal{A}}(\mathbf{r}) \in \arg \min_{\mathbf{x} \in \mathcal{A}} \langle \mathbf{r}, \mathbf{x} \rangle,$$

where $\mathcal{A} \subset \mathbb{R}^d$, $|\mathcal{A}| = m$ is a finite set of vectors that is usually referred as atoms. It is worth noticing that this linear minimization oracle (LMO) for FW step naturally chooses a single vector in \mathcal{A} that minimizes the inner product with \mathbf{x} . Again, this FW step creates some "partial update" that could be appreciated in many machine learning applications. Specifically, if our computation of gradient is again dominated by a matrix-vector (data matrix versus variable \mathbf{x}) inner product, we could possibly pre-compute each value of $\mathbf{v}_i := A\mathbf{x}_i$, $\mathbf{x}_i \in \mathcal{A}$, and simply use \mathbf{v}_i to update the gradient information when \mathbf{x}_i is the greedy direction provided by LMO.

When connecting to our sparse update case, we are now looking for a k -sparse update, $k \ll m = |\mathcal{A}|$, with the basis of \mathcal{A} , i.e., $\tilde{\mathbf{x}} = \sum_{i=1}^k \lambda_i \mathbf{x}_{n_i}$, $\mathbf{x}_{n_i} \in \mathcal{A}$. In this way, when we update $\mathbf{x}^+ \leftarrow (1 - \eta)\mathbf{x} + \eta\tilde{\mathbf{x}}$, we will only need to compute $\sum_{i=1}^k \mathbf{v}_{n_i}$ which is $\mathcal{O}(kd)$ time complexity.

However, to enforce such update that is "sparse" on \mathcal{A} is much harder. To migrate our algorithms with ℓ_1 ball or trace norm ball, we will essentially be solving the following problem:

$$\tilde{\mathbf{x}} \leftarrow \arg \min_{\Lambda \in \Delta^m, \|\Lambda\|_0 \leq k, \mathbf{x} = \sum_{i=1}^m \lambda_i \mathbf{x}_i, \mathbf{x}_i \in \mathcal{A}} \langle \mathbf{g}, \mathbf{y} \rangle + \frac{1}{2\eta} \|\mathbf{y} - \mathbf{x}\|_2^2,$$

where Δ^m is the m dimensional simplex, and \mathbf{g} is the current gradient vector.

Unlike the original sparse recovery problem that could be relaxed with an ℓ_1 constraint to softly encourage sparsity, it's generally much harder to find the k sparse Λ in this case. Actually, it is as hard as the lattice problem [19] and is NP hard in general.

Therefore we are not able to achieve linear convergence with cheap update with polytope-type constraints. Nonetheless, the naive FW with primal dual formulation should still be computational efficient in terms of per iteration cost, where a concentration on SVM on its dual form has been explored by [22].

B Discussions on Efficient Coordinate Selections

The modified Block Frank-Wolfe step in Eqn. (3) achieves an s -sparse update of the iterates and could be computed efficiently when one knows which s coordinates to update. However, in order to find the s coordinates, one needs to compute the full gradient $\nabla f(\mathbf{x})$ with naive implementation. This phenomenon reminds us of greedy coordinate descent.

Even with the known fact that coordinate descent converges faster with greedy selection than with random order[30], there have been hardness to propagate this idea because of expensive greedy selections since the arguments that GCD converges similarly with RCD in [28], except for special cases [25, 24, 6, 17]. This is also probability why the partial updates nature of FW steps is less exploited before.

We investigate some possible tricks to boost GCD method that could be possibly applied to FW methods. A recent paper [17], Karimireddy et al. make connections between the efficient choice of the greedy coordinates with the problem of Maximum Inner Product Search (MIPS) for a composite function $P(\mathbf{x}) = f(A\mathbf{x}) + g(\mathbf{x})$, where $A \in \mathbb{R}^{n \times d}$. We rephrase the connection for the Frank-Wolfe algorithm. Since the computation of gradient is essentially $A^\top \nabla f_{|A\mathbf{x}} + \nabla g(\mathbf{x})$, to find its largest magnitude is to search maximum inner products among:

$$\pm \langle [\tilde{\mathbf{a}}_i^\top | 1], [\nabla f_{|A\mathbf{x}}^\top | \nabla_i g(\mathbf{x})] \rangle, \text{ i.e. } \pm (\tilde{\mathbf{a}}_i^\top \nabla f_{|A\mathbf{x}} + \nabla_i g(\mathbf{x})),$$

where $\tilde{\mathbf{a}}_i \in \mathbb{R}^n$ is the i -th column of data matrix A , and $\nabla f_{|A\mathbf{x}}$ is the gradient of f at $A\mathbf{x}$. In this way, we are able to select the greedy coordinates by conducting MIPS for a fixed $\mathbb{R}^{2d \times (n+1)}$ matrix $[A^\top | I] - [A^\top | -I]^\top$ and each newly generated vector $[\nabla f_{|A\mathbf{x}}^\top | \nabla_i g(\mathbf{x})]$. Therefore when ∇g_i is constant for linear function or $\pm \lambda$ for $g(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$, we could find the largest magnitude of the gradient in sublinear time. Still, the problems it could conquer is very limited. It doesn't even work for ℓ_2 regularizer since the different coordinates in $\nabla_i g(\mathbf{x})$ creates d new vectors in each iteration and traditional MIPS could resolve it in time sublinear to d . Meanwhile, even with constant $\nabla_i g(\mathbf{x})$, it still requires at least $\mathcal{O}((2d)^c \log(d))$ times of inner products of dimension $n+1$ for some constant c [34].

However, we have shown that for general composite form $f(A\mathbf{x}) + g(\mathbf{x})$ with much more relaxed requirements on the regularizer g , we are able to select and update each coordinate with *constant* times of inner products on average while achieving linear convergence. Therefore the usage of these tricks applied on FW method (MIPS as well as the nearest neighbor search [6]) is completely dominated by our contribution and we omit them in the main text of this paper.

C More Results on Empirical Studies

C.1 More experiments with ℓ_1 norm

To investigate more on how our algorithms perform with different choices of parameters, we conducted more empirical studies with different settings of condition numbers. Specifically, we vary the parameter μ that controls the strong convexity of the primal function. Experiments are shown in Figure 2.

C.2 Experiments with trace norm ball on synthetic data

For trace norm constraints, we also implemented our proposal Primal Dual Block Frank Wolfe to compare with some prior work, especially Block FW [1]. Since prior work were mostly implemented in Matlab to tackle trace norm projections, we therefore also use Matlab to show fair comparisons. We choose quadratic loss $f(AX) = \|AX - B\|_F^2$ and g to be ℓ_2 regularizer with $\mu = 10/n$. The synthetic sensing matrix $A \in \mathbb{R}^{n \times d}$ is dense with $n = 1000$ and $d = 800$. Our observation B is of dimension 1000×600 and is generated by a ground truth matrix X_0 such that $B = AX_0$. Here $X_0 \in \mathbb{R}^{800 \times 600}$ is constructed with low rank structure. We vary its rank s to be 10, 20, and 100. The comparisons with stochastic FW, blockFW [1], STORC [13], SCGS [23], and projected SVRG [15] are presented in Figure 3, which verifies that our proposal PDBFW consistently outperforms the baseline algorithms.

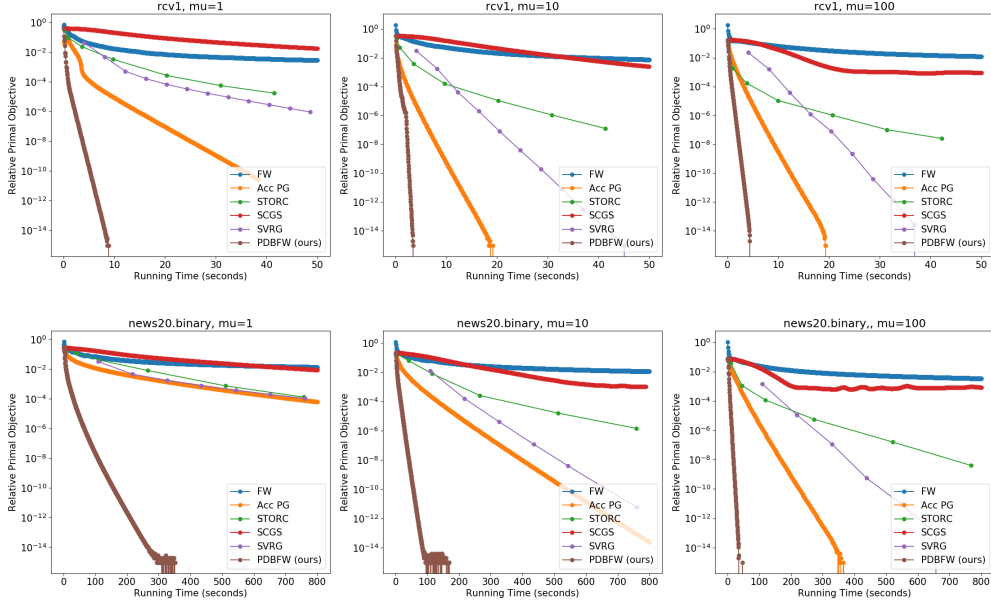


Figure 2: Convergence result comparison of different algorithms on smoothed hinge loss by varying the coefficient of the regularizer. The first row is the results ran on the rcv1.binary dataset, while the second row is the results ran on the news20.binary dataset. The first column is the result when the regularizer coefficient μ is set to $1/n$. The middle column is when $\mu = 10/n$, and the right column is when $\mu = 100/n$.

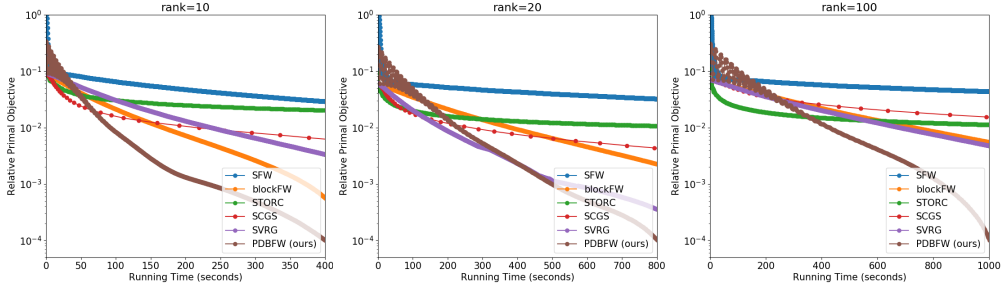


Figure 3: Convergence comparison of our Primal Dual Block Frank Wolfe and other baselines. Figures show the relative primal objective value decreases with the wall time.