

A. Appendix A: Convergence Analysis

A.1. Proof of Theorem 4.2

Recall primal, dual and Lagrangian forms:

$$P(\mathbf{x}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \phi_i(\langle A_i, \mathbf{x} \rangle) + g(\mathbf{x}), \quad (20)$$

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) \stackrel{\text{def}}{=} g(\mathbf{x}) + \frac{1}{n} \mathbf{y}^T A \mathbf{x} - \frac{1}{n} \sum_{i=1}^n \phi_i^*(y_i) \quad (21)$$

$$D(\mathbf{y}) \stackrel{\text{def}}{=} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}) \equiv \mathcal{L}(\bar{\mathbf{x}}(\mathbf{y}), \mathbf{y}) \quad (22)$$

where $\bar{\mathbf{x}}(\mathbf{y}) : \mathbb{R}^n \rightarrow \mathbb{R}^d$ is the optimal primal variable with respect to some \mathbf{y} , namely,

$$\bar{\mathbf{x}}(\mathbf{y}) = \arg \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})$$

For simplicity, we will use $\bar{\mathbf{x}}^{(t)} \stackrel{\text{def}}{=} \bar{\mathbf{x}}(\mathbf{y}^{(t)})$ throughout this paper. Similarly, we also use $\bar{\mathbf{y}}(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^n$ to be the optimal dual variable with respect to some \mathbf{x} .

Recall with the choice of regularizer of our model, $g(\mathbf{x}) = h(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$, where $h(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x}\|_2^2$ satisfies μ -strong convexity, μ -smooth and separable. The conjugate of loss function (e.g. smooth hinge loss used in our experiments): ϕ^* is γ -strongly convex.

Recall 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 $\Delta_d^{(t)} \stackrel{\text{def}}{=} D^* - D(\mathbf{y}^{(t)})$. In the proof, we will connect the objective change in primal/dual update with the primal/dual gap and show how the sub-optimality: $\Delta^{(t)} = \Delta_p^{(t)} + \Delta_d^{(t)}$ enjoys linear convergence.

Lemma A.1. (Primal Progress):

$$\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) \geq \frac{1}{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0 - 1} \Delta_p^{(t)}$$

Proof. This lemma is a direct result by our greedy update rule of our primal variables.

$$\begin{aligned} & \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)}) \\ &= \sum_i \{ \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}((\bar{x}_i^{(t)} - x_i^{(t)})\mathbf{e}_i + \mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \} \\ &= \sum_{i \in \text{supp}(\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)})} \{ \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}((\bar{x}_i^{(t)} - x_i^{(t)})\mathbf{e}_i + \mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \} \\ &\leq \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_0 \times \\ & \quad \max_i \{ \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}((\bar{x}_i^{(t)} - x_i^{(t)})\mathbf{e}_i + \mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \} \\ &= \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|_0 (\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)})) \end{aligned}$$

And by adding $\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t)}) - \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$ to both sides we finishes the proof. \square

Recall $i^{(t)}$ is the selected coordinate to update in dual variable $\mathbf{y}^{(t)}$.

Lemma A.2. (Primal-Dual Progress).

$$\begin{aligned} & \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\ &\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\ & \quad + \eta \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle_{i^{(t)}}^2 - \eta \left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g \right)^2 \right) \end{aligned}$$

, where $g \in \frac{1}{n} \partial \phi_{i^{(t)}}^*(\mathbf{y}^{(t)})$.

Our goal is to prove that $\Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \leq -\delta \Delta_p^{(t)} - \delta \Delta_d^{(t)}$ to show linear convergence in sub-optimality. Since $\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \leq -\frac{1}{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0} \Delta_p^{(t)}$, this lemma is the middle step that connects to the primal part, and the remaining part represents the dual progress and will be analyzed later.

Proof. The primal and dual gap comes from both primal and dual progresses:

$$\underbrace{\Delta_d^{(t)} - \Delta_d^{(t-1)}}_{\text{dual progress}} + \underbrace{\Delta_p^{(t)} - \Delta_p^{(t-1)}}_{\text{primal progress}}$$

- Dual progress:

By Danskins' theorem, $-D(\mathbf{y})$ is γ -strongly convex. Therefore for any $g \in \partial \phi_{i^{(t)}}^*(\mathbf{y}^{(t)})$, we have,

$$\begin{aligned} & \Delta_d^{(t)} - \Delta_d^{(t-1)} = (-D(\mathbf{y}^{(t)}) - (-D(\mathbf{y}^{(t-1)}))) \\ & \leq -\left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g \right) (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)}) \\ & \quad - \frac{\gamma}{2} (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)})^2 \end{aligned} \quad (23)$$

- Primal progress:

Similarly we get,

$$\begin{aligned} & \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^{(t-1)}) \\ & \leq \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - g \right) (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)}) \\ & \quad + \frac{\gamma}{2} (y_{i^{(t)}}^{(t-1)} - y_{i^{(t)}}^{(t)})^2 \end{aligned} \quad (24)$$

Therefore,

$$\begin{aligned} & \Delta_p^{(t)} - \Delta_p^{(t-1)} \\ &= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^{(t-1)}) - (D(\mathbf{y}^{(t)}) - D(\mathbf{y}^{(t-1)})) \\ &= \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) + \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^{(t-1)}) \\ & \quad - (D(\mathbf{y}^{(t)}) - D(\mathbf{y}^{(t-1)})) \\ &\leq \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\ & \quad + \frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)}) \end{aligned}$$

Here the last inequality comes from inequalities (24) and (23).

Meanwhile, with the update rule of dual variable:

$$y_{i^{(t)}}^{(t)} \leftarrow \arg \max_{\beta} \frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle \beta - \phi_{i^{(t)}}^*(\beta) - \frac{1}{2\eta} (\beta - y_{i^{(t)}}^{(t)})^2$$

Therefore $\exists g \in \partial \phi_{i^{(t)}}^*(\mathbf{y}^{(t)})$ such that $y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)} = \eta \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - g \right)$. Therefore:

$$\begin{aligned} (23) &= -\left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g \right) (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)}) \\ & \quad - \frac{\gamma}{2} (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)})^2 \\ &= \left\langle \frac{1}{n} A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)} \right\rangle (y_{i^{(t)}}^{(t-1)} - y_{i^{(t)}}^{(t)}) \\ & \quad - \left(\frac{1}{\eta} + \frac{\gamma}{2} \right) (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)})^2 \end{aligned} \quad (25)$$

- Summing together we have:

$$\begin{aligned}
 & \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
 \leq & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\
 & + \frac{2}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)} \rangle (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)}) \\
 & - \left(\frac{1}{\eta} + \frac{\gamma}{2}\right) (y_{i^{(t)}}^{(t)} - y_{i^{(t)}}^{(t-1)})^2 \\
 = & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\
 & + \frac{2\eta}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - g\right) \\
 & - \eta^2 \left(\frac{1}{\eta} + \frac{\gamma}{2}\right) \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - g\right)^2 \\
 \leq & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\
 & + \eta \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle\right)^2 - \eta \left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g\right)^2
 \end{aligned}$$

□

Afterwards, we upper bound the dual progress $(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle)^2 - (\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g)^2$ by dual gap $\Delta_d^{(t)}$:

Lemma A.3. (Dual Progress).

$$\begin{aligned}
 & \left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)} \rangle\right)^2 - \left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - g\right)^2 \\
 \leq & -\frac{\gamma}{2n} \Delta_d^{(t)} + \frac{5R^2}{2n^2} \|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|^2
 \end{aligned} \tag{26}$$

, where $g \in \frac{1}{n} \partial \phi_{i^{(t)}}^*(y_{i^{(t)}}^{(t)})$.

Proof. For simplicity, we denote $\phi^*(\mathbf{y}) = \frac{1}{n} \sum_i \phi_i^*(y_i)$. To begin with,

$$\begin{aligned}
 \Delta_d^{(t)} = D^* - D(\mathbf{y}) & \leq \frac{2}{\gamma} \left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|^2 \\
 & \leq \frac{2n}{\gamma} \left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_\infty^2
 \end{aligned}$$

In our algorithm, the greedy choice of $i^{(t)}$ makes sure $\left\| \frac{1}{n} A \mathbf{x}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_{i^{(t)}} = \left\| \frac{1}{n} A \mathbf{x}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_\infty$. However, here we need the relation between $\left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_{i^{(t)}}$ and $\left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_\infty$ (assumed to be reached at coordinate i^*). We bridge their gap by $\delta \stackrel{\text{def}}{=} \frac{1}{n} A(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)})$. Since

$$\begin{aligned}
 & -\left(\frac{1}{n} \langle A_{i^{(t)}}, \bar{\mathbf{x}}^{(t)} \rangle - \frac{1}{n} (\phi_{i^{(t)}}^*)'(y_{i^{(t)}}^{(t)})\right)^2 \\
 = & -\left(\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - \frac{1}{n} (\phi_{i^{(t)}}^*)'(y_{i^{(t)}}^{(t)}) + \delta_{i^{(t)}}\right)^2 \\
 \leq & -\frac{1}{2n^2} \left(\langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - (\phi_{i^{(t)}}^*)'(y_{i^{(t)}}^{(t)}) \right)^2 + \delta_{i^{(t)}}^2 \\
 = & -\frac{1}{2} \left\| \frac{1}{n} A \mathbf{x}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_\infty^2 + \delta_{i^{(t)}}^2 \\
 \leq & -\frac{1}{2} \left(\frac{1}{n} \langle A_{i^*}, \mathbf{x}^{(t)} \rangle - \frac{1}{n} (\phi_{i^*}^*)'(y_{i^*}^{(t)}) \right)^2 + \|\delta\|_\infty^2 \\
 = & -\frac{1}{2} \left(\frac{1}{n} \langle A_{i^*}, \bar{\mathbf{x}}^{(t)} \rangle - \frac{1}{n} (\phi_{i^*}^*)'(y_{i^*}^{(t)}) - \delta_{i^*} \right)^2 + \|\delta\|_\infty^2 \\
 \leq & -\frac{1}{4} \left(\frac{1}{n} \langle A_{i^*}, \bar{\mathbf{x}}^{(t)} \rangle - \frac{1}{n} (\phi_{i^*}^*)'(y_{i^*}^{(t)}) \right)^2 + \frac{3}{2} \|\delta\|_\infty^2 \\
 = & -\frac{1}{4} \left\| \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \partial \phi^*(\mathbf{y}^{(t)}) \right\|_\infty^2 + \frac{3}{2} \|\delta\|_\infty^2 \\
 \leq & -\frac{\gamma}{2n} \Delta_d^{(t)} + \frac{3}{2} \|\delta\|_\infty^2
 \end{aligned}$$

The first inequality follows $-(a+b)^2 = -a^2 - b^2 - 2ab \leq -a^2 - b^2 + \frac{1}{2}a^2 + 2b^2 = -\frac{1}{2}a^2 + b^2$, and replace a by $\frac{1}{n} \langle A_{i^{(t)}}, \mathbf{x}^{(t)} \rangle - \frac{1}{n} (\phi_{i^{(t)}}^*)'(y_{i^{(t)}}^{(t)})$ and $b \stackrel{\text{def}}{=} \delta_{i^{(t)}}$. And similarly for the third inequality.

Meanwhile, since $\|A(\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)})\|_\infty \leq R \|\bar{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)}\|$, together we get Lemma A.3. □

Now we have established the connection between the primal and dual progress (change in primal/dual gap) with primal and dual gap, and the only redundant part is $\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|$, but since $\frac{\mu}{2} \|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\| \leq \mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})$, which could be absorbed in the primal gap. Therefore, back to the main inequality (26):

Proof of Theorem 4.2.

$$\begin{aligned}
 & \Delta_d^{(t)} - \Delta_d^{(t-1)} + \Delta_p^{(t)} - \Delta_p^{(t-1)} \\
 \stackrel{\text{Lemma A.2}}{\leq} & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) \\
 & + \left\langle \frac{1}{n} A \mathbf{x}^{(t)} - \nabla \varphi(\mathbf{y}^{(t)}), \mathbf{y}^{(t)} - \mathbf{y}^{(t-1)} \right\rangle \\
 & - 2 \left\langle \frac{1}{n} A \bar{\mathbf{x}}^{(t)} - \nabla \varphi(\mathbf{y}^{(t)}), \mathbf{y}^{(t)} - \mathbf{y}^{(t-1)} \right\rangle \\
 \stackrel{\text{Lemma A.3}}{\leq} & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) - \frac{\eta\gamma}{2n} \Delta_d^{(t)} \\
 & + \frac{5\eta R^2}{2n^2} \|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|^2 \\
 \leq & \mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t) - \frac{\eta\gamma}{2n} \Delta_d^{(t)} \\
 & + \frac{5\eta R^2}{\mu n^2} (\mathcal{L}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^{(t)})) \\
 = & \left(1 - \frac{5\eta R^2}{\mu n^2}\right) (\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{y}^t)) - \frac{\eta\gamma}{2n} \Delta_d^{(t)} \\
 & + \frac{5\eta R^2}{\mu n^2} (\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^t)) \\
 \stackrel{\text{Lemma A.1}}{\leq} & -\left(1 - \frac{5\eta R^2}{\mu n^2}\right) \frac{1}{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0 - 1} \Delta_p^{(t)} - \frac{\eta\gamma}{2n} \Delta_d^{(t)} + \\
 & \frac{5\eta R^2}{\mu n^2} (\mathcal{L}(\mathbf{x}^{(t+1)}, \mathbf{y}^t) - \mathcal{L}(\bar{\mathbf{x}}^{(t)}, \mathbf{y}^t)) \\
 = & -\left(\left(1 - \frac{5\eta R^2}{\mu n^2}\right) \frac{1}{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0 - 1} - \frac{5\eta R^2}{\mu n^2}\right) \Delta_p^{(t)} \\
 & - \frac{\eta\gamma}{2n} \Delta_d^{(t)}
 \end{aligned}$$

Therefore, we have

$$\frac{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0}{\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0 - 1} \left(1 - \frac{5\eta R^2}{\mu n^2}\right) \Delta_p^{(t)} + \left(1 + \frac{\eta\gamma}{2n}\right) \Delta_d^{(t)} \leq \Delta_d^{(t-1)} + \Delta_p^{(t-1)}$$

i.e. linear convergence. Notice when

$$\begin{aligned}
 \eta^{(t)} & \leq \frac{2n^2 \mu}{(10R^2 + n\gamma\mu) \|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0} \\
 \Delta^{(t)} & \leq \frac{1}{1 + \frac{\eta^{(t)}\gamma}{2n}} \Delta^{(t-1)}
 \end{aligned} \tag{27}$$

Specifically, when inequality holds for (27), and suppose $\|\mathbf{x}^{(t)} - \bar{\mathbf{x}}^{(t)}\|_0 \leq s$, then it requires $\mathcal{O}(s(\frac{\kappa}{n} + 1) \log \frac{1}{\epsilon})$ iterations to achieve ϵ primal and dual sub-optimality, where $\kappa = \frac{R^2}{\mu\gamma}$. □

B Appendix B: Additional Experimental Results

Finally, we show result for $\lambda = 0.01, 0.1$, and $\mu = 0.01, 0.1, 1$. Here are some comments for results under different parameters.

The winning margin of DGPD is larger on data sets of dense feature matrix than that of sparse feature matrix. One reason for this is, for data of sparse feature matrix, features of higher frequency are more likely to be active than those of lower frequency, and therefore, the feature sub-matrix corresponding to the *active primal variables* are often denser than submatrix corresponding to the inactive ones. This results in a less overall speedup.

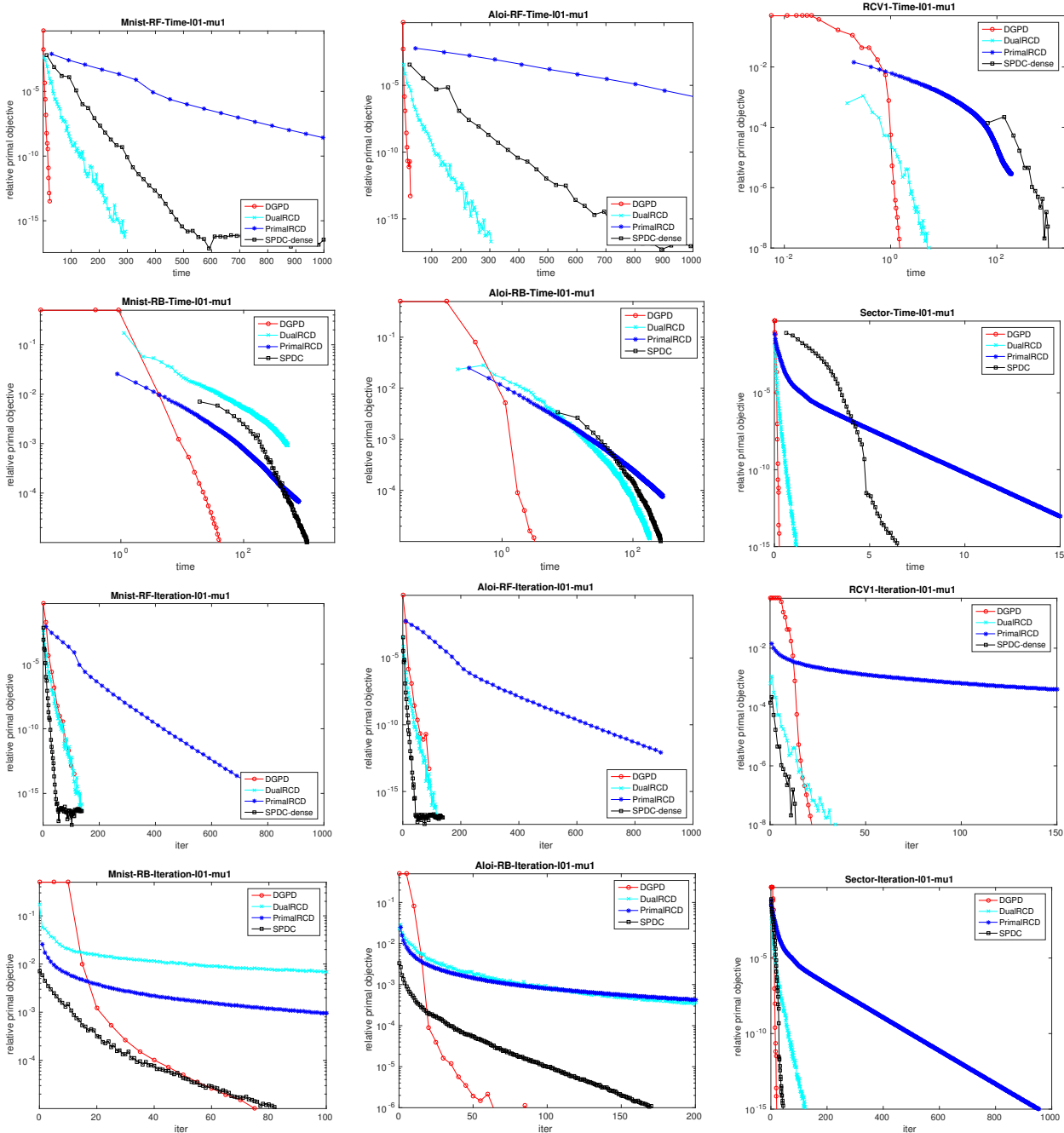


Figure 2. Relative Objective versus Time (the upper 2 rows) and versus # iterations (the lower 2 rows) for $\lambda = 0.1, \mu = 1$.

We also observe that in order to achieve the best performance of DGPD, both primal and dual sparsity must hold, and the sparsity is partially controlled by the L1/L2 penalty. In particular, when the L1 penalty has too much weight, the primal iterate would become too sparse to yield a reasonable prediction accuracy, which then results in a particularly dense dual iterate due to its non-zero loss on most of the samples. Another example is, when the L2 penalty becomes too large, the classifier would tend to mis-classify many examples in order to gain a large margin, which results in dense dual iterates.

However, in practice such hyperparameter settings are less likely to be chosen due to its inferior prediction performance.

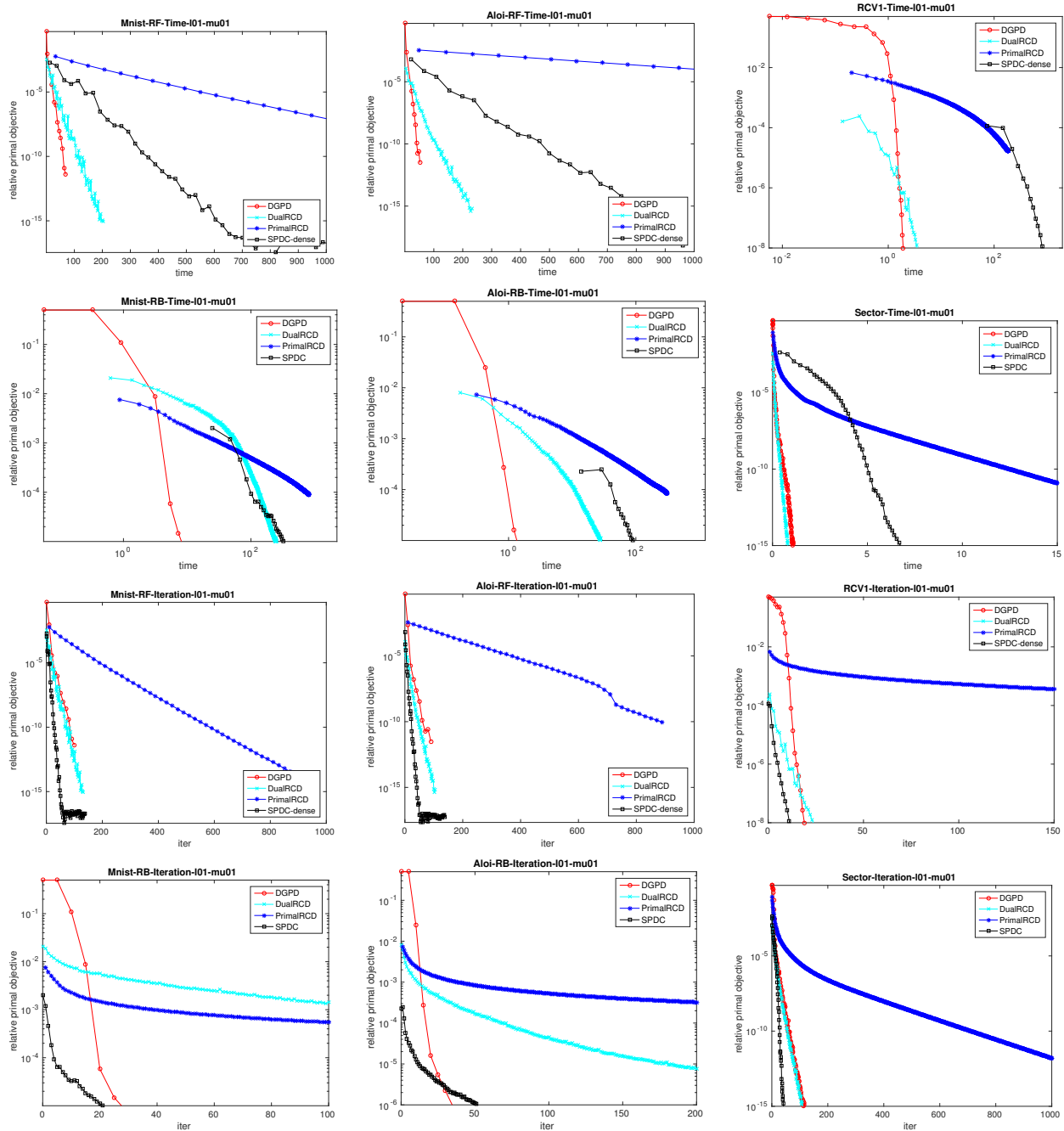


Figure 3. Relative Objective versus Time (the upper 2 rows) and versus # iterations (the lower 2 rows) for $\lambda = 0.1, \mu = 0.1$.

Doubly Greedy Primal-dual Coordinate Descent for Sparse Empirical Risk Minimization

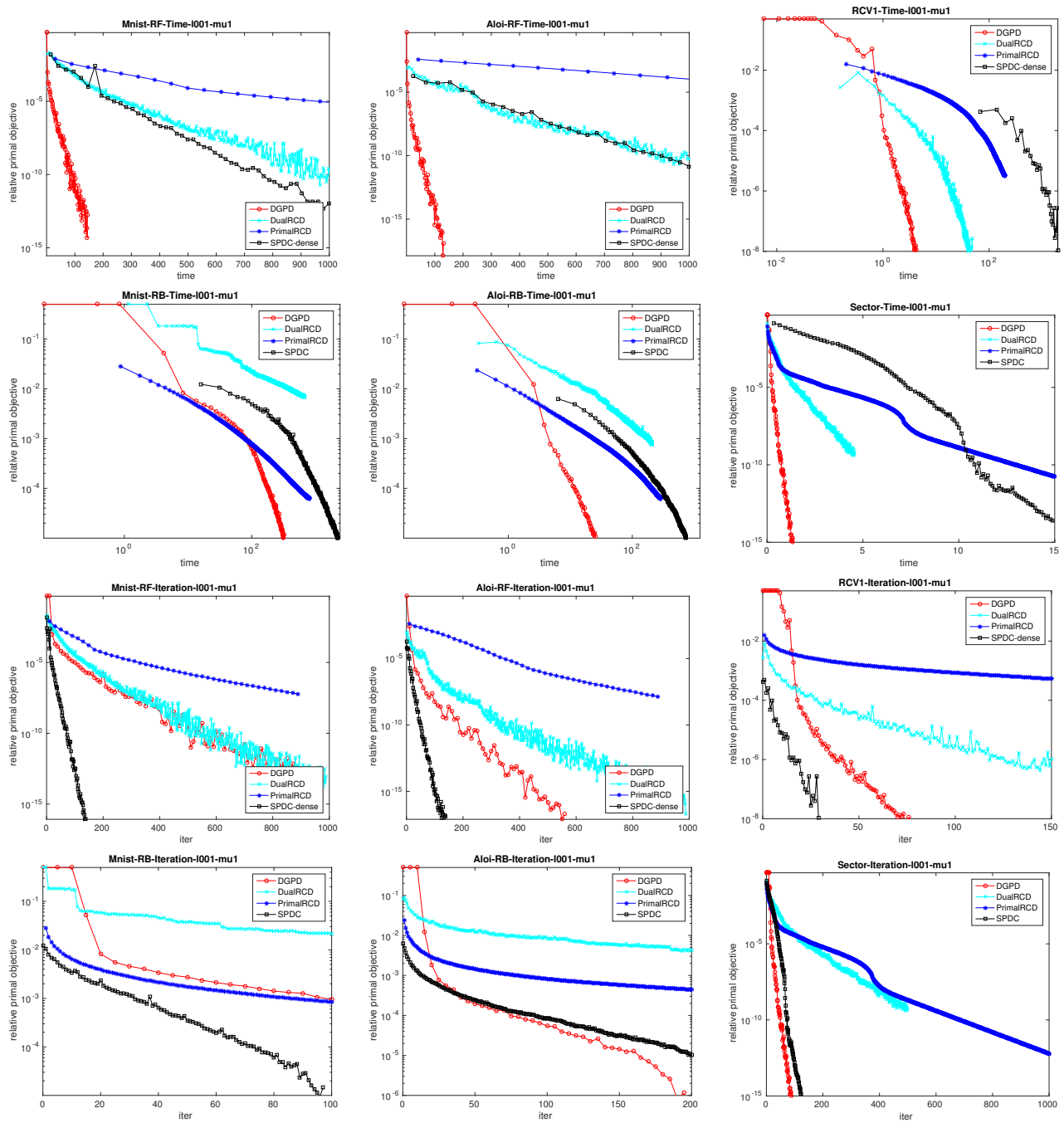


Figure 4. Relative Objective versus Time (the upper 2 rows) and versus # iterations (the lower 2 rows) for $\lambda = 0.01$, $\mu = 1$.

Doubly Greedy Primal-dual Coordinate Descent for Sparse Empirical Risk Minimization

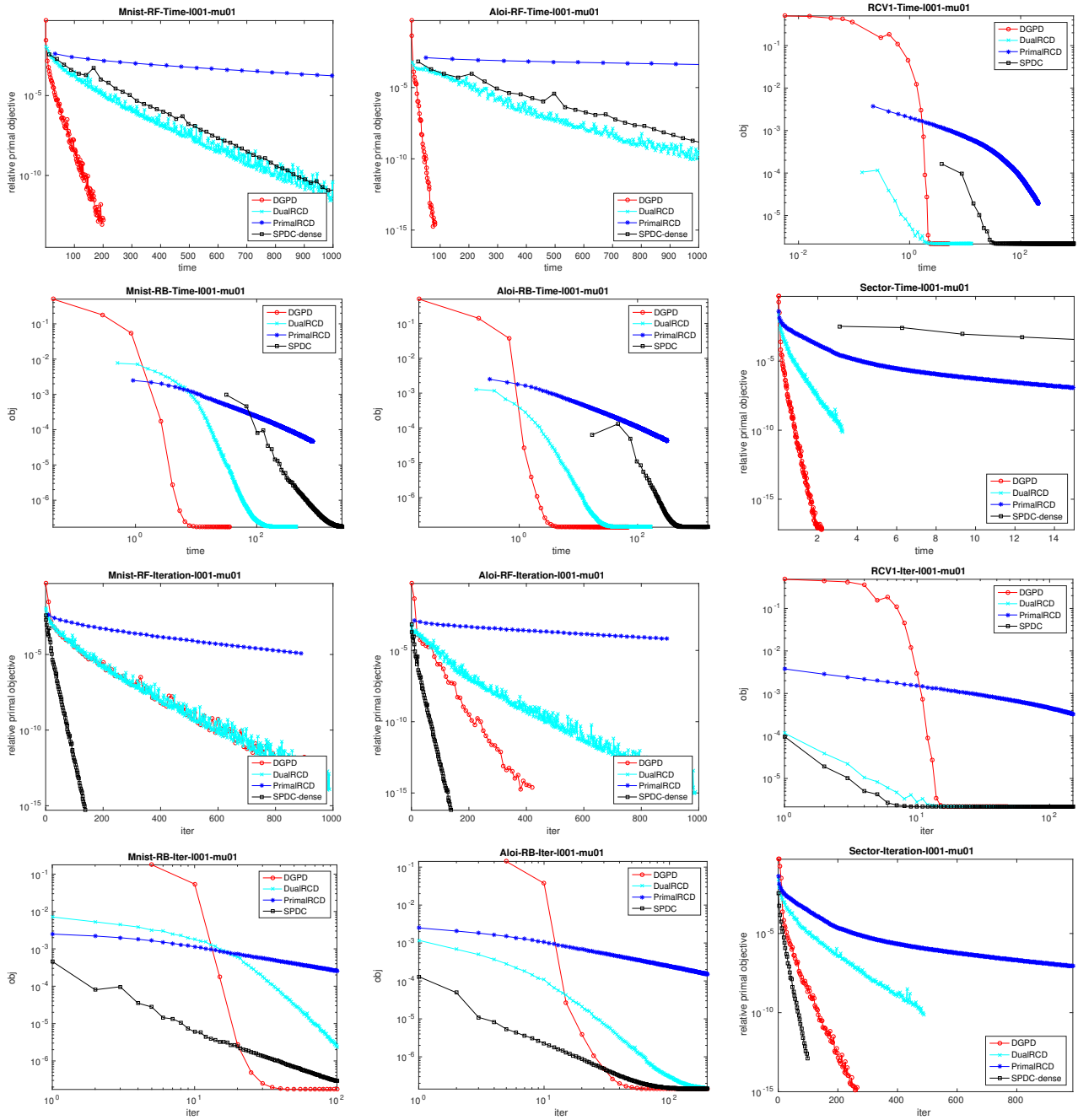


Figure 5. Relative Objective versus Time (the upper 2 rows) and versus # iterations (the lower 2 rows) for $\lambda = 0.01$, $\mu = 0.1$.

Doubly Greedy Primal-dual Coordinate Descent for Sparse Empirical Risk Minimization

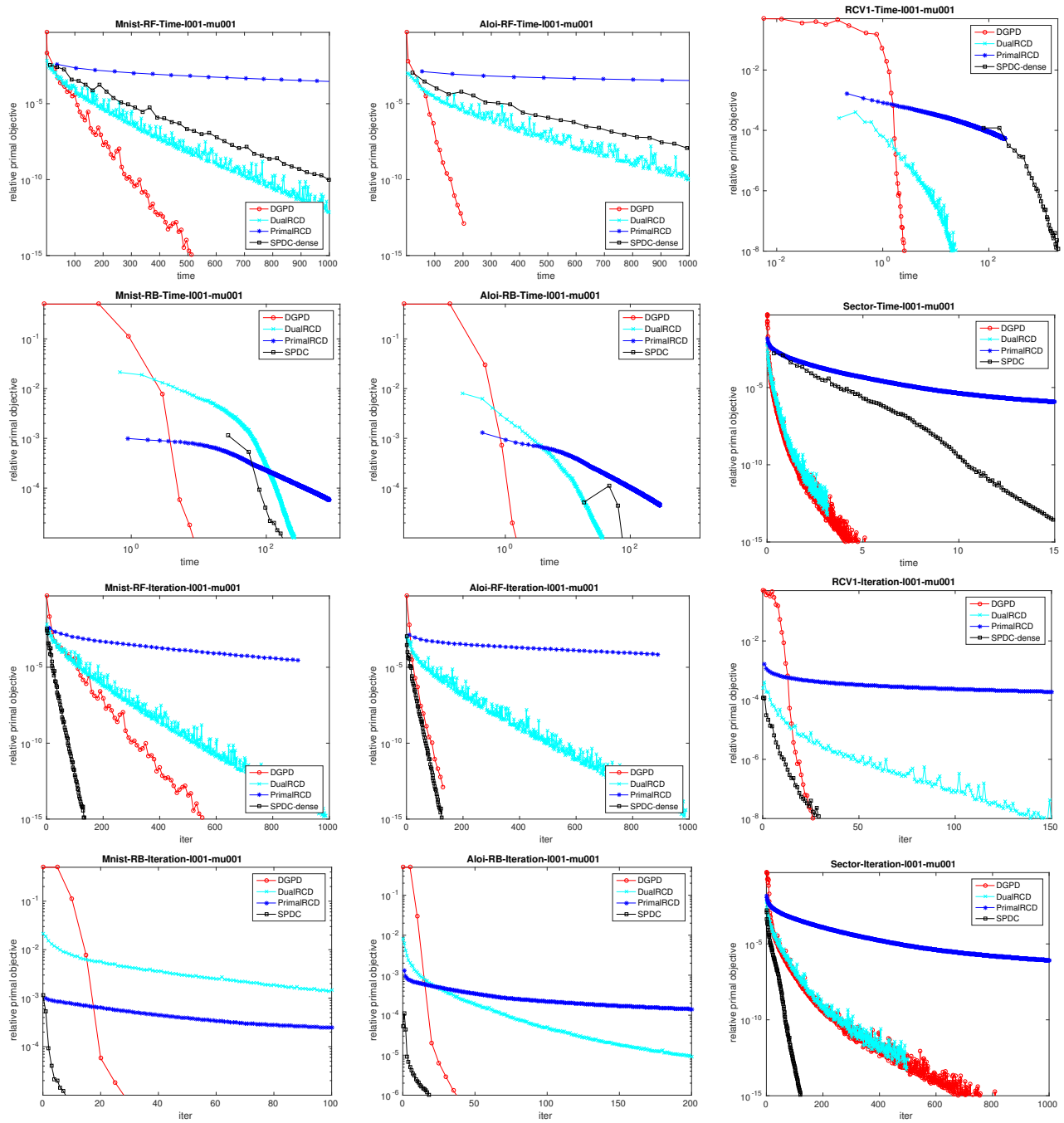


Figure 6. Relative Objective versus Time (the upper 2 rows) and versus # iterations (the lower 2 rows) for $\lambda = 0.01$, $\mu = 0.01$.