

## 7 Appendix

### 7.1 Proof of Theorem 1

*Proof.* There are two main difficulties in proving the convergence of our algorithm, and none of them is addressed in previous works. First, the Hessian matrix  $\mathcal{H}$  is a block-structured matrix as shown in (7), and unfortunately it is low-rank. Second, we have to show the convergence with the active set selection technique.

Let  $H(\boldsymbol{\theta}) \in \mathcal{R}^{n \times n}$  be the Hessian  $\nabla^2 \mathcal{L}(\bar{\boldsymbol{\theta}})$  when  $\mathcal{L}(\cdot)$  is strongly convex. When it is not, as discussed in Section 3, we set  $H = \nabla^2 \mathcal{L}(\bar{\boldsymbol{\theta}}) + \epsilon I$  for a small constant  $\epsilon > 0$ . Let  $\mathbf{d} \in \mathbb{R}^{n^2 k}$  denotes the minimizer of the quadratic subproblem (8). By definition, we can easily observe that  $\sum_{r=1}^k \mathbf{d}^{(r)} = \mathbf{y}$  where

$$\mathbf{y}^* := \operatorname{argmin}_{\mathbf{y}} \mathbf{y}^T H(\boldsymbol{\theta}) \mathbf{y} + \langle \mathbf{y}, \nabla \mathcal{L}(\bar{\boldsymbol{\theta}}) \rangle + \mathcal{W}(\mathbf{y}), \quad (19)$$

where  $\mathcal{W}(\mathbf{y}) = \min_{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(k)}: \sum_{r=1}^k \mathbf{y}^{(r)} = \mathbf{y}} \sum_{r=1}^k \lambda_r \|\mathbf{y}^{(r)}\|_{\mathcal{A}_r}$ . Now (19) has a strongly convex Hessian, and thus if  $\mathcal{W}(\mathbf{y})$  is convex then Theorem 3.1 in [16] can be used to show that the algorithm converges, and thus  $F(\sum_{i=1}^r \boldsymbol{\theta}^{(i)})$  converges (without active subspace selection).

To show  $\mathcal{W}$  is convex, assume we have  $\mathbf{a}, \mathbf{b}$  and  $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(k)}$  are decomposition of  $\mathbf{a}$  that attains the minimizer of  $\mathcal{W}(\mathbf{a})$ . By definition we have

$$\begin{aligned} \mathcal{W}(\alpha \mathbf{a} + (1 - \alpha) \mathbf{b}) &\leq \sum_{r=1}^k \lambda_r \|\alpha \mathbf{a}^{(r)} + \mathbf{b}^{(r)}\|_{\mathcal{A}_r} \\ &\leq \sum_{r=1}^k \lambda_r (\alpha \|\mathbf{a}^{(r)}\|_{\mathcal{A}_r} + (1 - \alpha) \|\mathbf{b}^{(r)}\|_{\mathcal{A}_r}) \\ &= \alpha \mathcal{W}(\mathbf{a}) + (1 - \alpha) \mathcal{W}(\mathbf{b}). \end{aligned}$$

Thus  $\mathcal{W}$  is convex, and if we solve each quadratic approximation exactly without active subspace selection, our algorithm converges to the global optimum.

Next we discuss the convergence of our algorithm with active subspace selection. [12] has shown the convergence of active set selection under  $\ell_1$  regularization, but here the situation is different – we can have infinite many atoms, as in group lasso or nuclear norm cases, so the original proof cannot be applied. To analyze the active subspace selection technique, we will use the convergence proof for Block Coordinate Gradient Descent (BCGD) in [23]. To begin, we give a quick review of the Block Coordinate Gradient Descent (BCGD) algorithm discussed in [23].

BCGD is proposed to solve the composite functions with the following form:

$$\operatorname{argmin}_{\mathbf{x}} F(\mathbf{x}) := f(\mathbf{x}) + P(\mathbf{x}),$$

where  $P(\mathbf{x})$  is a convex and separable function. Here we consider  $\mathbf{x}$  to be a  $p$  dimensional vector, but in general  $\mathbf{x}$  can be in any space. At the  $t$ -th iteration, BCGD chooses a subset  $\mathcal{J}_t$  and compute the descent direction by

$$\mathbf{d}_t = \mathbf{d}_{H_t}^{\mathcal{J}_t}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{d}} \left\{ \langle \nabla f(\mathbf{x}), \mathbf{d} \rangle + \frac{1}{2} \mathbf{d}^T H_t \mathbf{d} + P(\mathbf{x} + \mathbf{d}) \mid d_i = 0, \forall i \notin \mathcal{J}_t \right\} \equiv \operatorname{argmin}_{\mathbf{d}} Q_{H_t}^{\mathcal{J}_t}(\mathbf{d}). \quad (20)$$

After computing the descent direction  $\mathbf{d}_{H_t}^{\mathcal{J}_t}$ , a line search procedure is used to find the descent direction, where the largest  $\alpha_t$  is chosen by searching over  $1, 1/2, 1/4, \dots$  satisfying

$$F(\mathbf{x}_t + \alpha \mathbf{d}_t) \leq F(\mathbf{x}_t) + \alpha_t \sigma \Delta_t, \text{ where } \Delta_t = \langle \nabla f(\mathbf{x}_t), \mathbf{d}_t \rangle + P(\mathbf{x}_t + \mathbf{d}_t) - P(\mathbf{x}_t),$$

and  $\sigma \in (0, 1)$  is any constant. The line search condition is exactly the same with ours in (10).

It is shown in Theorem 3.1 of [23] that when  $\mathcal{J}_t$  is selected in a cyclic order covering all the indexes, than BCGD converges to the global optimum for any convex function  $f(X)$ .

To proof this theorem, we want to show the equivalence between our algorithm and BCGD with one block (BCGD-1block). We first prove the convergence for the case where we only have one regularizer, thus the problem is

$$\operatorname{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}) + \lambda \|\mathbf{x}\|_{\mathcal{A}}.$$

The key idea is to carefully define the matrix  $H_t$  at each iteration, according to our subspace selection trick. Note that in (20) the matrix  $H_t$  can be any positive definite matrix instead of Hessian. At each iteration of Algorithm 1, assume the fixed and free subspace defined in (13) are  $\mathcal{S}_{\text{fixed}}$  and  $\mathcal{S}_{\text{free}}$ . Assume  $\dim(\mathcal{S}_{\text{free}}) = q$  and  $\dim(\mathcal{S}_{\text{fixed}}) = n - q$ ,  $R = [\mathbf{r}_1, \dots, \mathbf{r}_q]$  are the orthogonal basis for  $\mathcal{S}_{\text{free}}$ , then we can construct  $\tilde{H}_t$  by

$$\tilde{H}_t = R(R^T H_t R)^T + R^\perp (R^\perp)^\perp. \quad (21)$$

It is easy to see that  $\tilde{H}_t$  is positive definite if the real Hessian  $\nabla^2 f(\mathbf{x})$  is positive definite. Using  $\tilde{H}_t$  as the Hessian in (20), we first consider  $\mathbf{d}_{\text{free}}$  to be the minimizer of the following problem:

$$\mathbf{d}_{\text{free}} = \operatorname{argmin}_{\Delta_{\text{free}} \in \mathcal{S}_{\text{free}}} \left\{ \langle \nabla f(\mathbf{x}), \Delta_{\text{free}} \rangle + \frac{1}{2} \Delta_{\text{free}}^T \tilde{H}_t \Delta_{\text{free}} + P(\mathbf{x} + \Delta_{\text{free}}) \right\},$$

which is the quadratic subproblem (20) within the subspace  $\mathcal{S}_{\text{free}}$ . Next we show that  $\mathbf{d}_{\text{free}}$  is indeed the optimizer of the whole quadratic subproblem (20) with the original Hessian  $H_t$ . To show this, taking derivative of (20) with respect to a  $\mathbf{a} \in \mathcal{S}_{\text{fixed}}$ , the subgradient will be

$$\partial_{\mathbf{a}} Q_{\tilde{H}_t}(\mathbf{d}_{\text{free}}) = \langle \nabla \mathcal{L}(\mathbf{x}), \mathbf{a} \rangle + \mathbf{a}^T \tilde{H}_t \mathbf{d}_{\text{free}} + \partial_{\mathbf{a}} P(\mathbf{x} + \mathbf{d}_{\text{free}}).$$

By the definition of  $\tilde{H}_t$  in (21), since  $\mathbf{a} \in \operatorname{span}(R^\perp)$  and  $\mathbf{d}_{\text{free}} \in \operatorname{span}(R)$ , we have  $\mathbf{a}^T \tilde{H}_t \mathbf{d}_{\text{free}} = 0$ . Also, since  $\mathbf{a} \in \mathcal{S}_{\text{fixed}}$ ,  $\langle \mathbf{x} + \mathbf{d}_{\text{free}}, \mathbf{a} \rangle = 0$ , thus  $\partial_{\mathbf{a}} P(\mathbf{x} + \mathbf{d}_{\text{free}}) = [-\lambda, \lambda]$ . Also, since  $\mathbf{a} \in \mathcal{S}_{\text{fixed}}$ ,  $|\langle \nabla f(\mathbf{x}), \mathbf{a} \rangle| < \lambda$ . Therefore, we have

$$0 \in \partial_{\mathbf{a}} Q_{\tilde{H}_t}(\mathbf{d}_{\text{free}}), \quad \forall \mathbf{a} \in \mathcal{S}_{\text{fixed}},$$

also, since  $\mathbf{d}_{\text{free}}$  is the optimal solution in  $\mathcal{S}_{\text{free}}$ , the projection of subgradient to  $\mathbf{a} \in \mathcal{S}_{\text{free}}$  is 0. Therefore  $\mathbf{d}_{\text{free}}$  is the minimizer of  $Q_{\tilde{H}_t}$ .

Based on the above discussion, our algorithm that computing the generalized Newton direction in free subspace  $\mathcal{S}_{\text{free}}$  is equivalent to another BCGD-1block algorithm with  $\tilde{H}_t$  as the approximated Hessian matrix. Therefore based on Theorem 3.1 of [23], our algorithm converges to the global optimum.

When there are more than one set of parameters, i.e., we want to solve (1) with parameter sets  $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(k)}$ . In this case, the Hessian matrix is a  $kn$  by  $kn$  matrix. For each parameter set, we will select  $\mathcal{S}_{\text{free}}^{(i)}$  and  $\mathcal{S}_{\text{fixed}}^{(i)}$ , and only update on  $\mathcal{S}_{\text{free}}^{(i)}$ . To prove the convergence, similar to the above arguments, we can construct the approximate Hessian  $\tilde{H}_t$  and show the equivalence between BCGD-1block with  $\tilde{H}_t$  and our algorithm.  $\tilde{H}_t$  can be divided into  $k^2$  blocks, each one is a  $n$  by  $n$  matrix  $\tilde{H}_t^{(i,j)}$ , where

$$\tilde{H}_t^{(i,j)} = R^{(i)} (R^{(i)})^T H R^{(j)} (R^{(j)})^T + (R^{(i)})^\perp ((R^{(j)})^\perp)^T,$$

where  $H = \nabla^2 \mathcal{L}(\sum_{i=1}^k \boldsymbol{\theta}^{(r)})$ ,  $R^{(i)}$  is the basis of  $\mathcal{S}_{\text{free}}^{(i)}$ , and  $(R^{(i)})^\perp$  is the basis of  $\mathcal{S}_{\text{fixed}}^{(i)}$ . Assume  $\mathbf{d}_{\text{free}} \in \mathcal{R}^{nk}$  is the solution of (20) with the constraint that projection to  $\mathcal{S}_{\text{fixed}}^{(i)}$  is 0 for all  $i$ , then it is the solution for both BCGD-1block with  $\tilde{H}_t$  and also the update direction produced by Algorithm 1. Also,

$$\mathbf{a} \tilde{H}_t \mathbf{d}_{\text{free}} = 0 \quad \forall \mathbf{a} \in \{\mathcal{S}_{\text{fixed}}^{(1)}, \dots, \mathcal{S}_{\text{fixed}}^{(k)}\},$$

therefore similar to the previous arguments we can show that

$$0 \in \partial_{\mathbf{a}} Q_{\tilde{H}_t}(\mathbf{d}_{\text{free}}), \quad \forall \mathbf{a} \in \{\mathcal{S}_{\text{fixed}}^{(1)}, \dots, \mathcal{S}_{\text{fixed}}^{(k)}\},$$

which indicates  $\mathbf{d}_{\text{free}}$  is the minimizer of  $Q_{\tilde{H}_t}(\mathbf{d})$ . Then we can see that the BCGD-1block algorithm with  $\tilde{H}_t$  as the Hessian matrix is equivalent to Algorithm 1 when each quadratic subproblem is solved exactly, therefore our algorithm converges to the global optimum of (1)  $\square$

## 7.2 Proof of Lemma 2

*Proof.* First, since  $\mathcal{L}(\cdot)$  is strongly convex,

$$\langle \nabla \mathcal{L}(\bar{\mathbf{x}}) - \nabla \mathcal{L}(\bar{\mathbf{y}}), \bar{\mathbf{x}} - \bar{\mathbf{y}} \rangle \geq \eta \|\bar{\mathbf{x}} - \bar{\mathbf{y}}\|^2. \quad (22)$$

Next, by the optimal condition we know

$$\begin{aligned} -\nabla\mathcal{L}(\bar{\mathbf{x}}) &\in \partial\|\mathbf{x}^{(r)}\|_{\mathcal{A}_r}, \forall r = 1, \dots, k \\ -\nabla\mathcal{L}(\bar{\mathbf{y}}) &\in \partial\|\mathbf{y}^{(r)}\|_{\mathcal{A}_r}, \forall r = 1, \dots, k. \end{aligned}$$

By the convexity for each  $\|\cdot\|_{\mathcal{A}_r}$ ,

$$\langle -\nabla\mathcal{L}(\bar{\mathbf{x}}) + \nabla\mathcal{L}(\bar{\mathbf{y}}), \mathbf{x}^{(r)} - \mathbf{y}^{(r)} \rangle \geq 0, \forall r = 1, \dots, k.$$

Summing  $r$  from 1 to  $k$  we get  $\langle -\nabla\mathcal{L}(\bar{\mathbf{x}}) + \nabla\mathcal{L}(\bar{\mathbf{y}}), \bar{\mathbf{x}} - \bar{\mathbf{y}} \rangle \geq 0$ , and combined with (22) we have  $\|\bar{\mathbf{x}} - \bar{\mathbf{y}}\| = 0$ .  $\square$

### 7.3 Proof of Theorem 2

*Proof.* Since assumption (16) holds, there exists a positive constant  $\epsilon$  such that

$$\|\Pi_{(\mathcal{T}^{(r)})^\perp}(\nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*))\|_{\mathcal{A}_r}^* < \lambda_r - \epsilon \quad \forall r = 1, \dots, k. \quad (23)$$

We first focus on showing that  $\mathcal{S}_{\text{fixed}}$  will be eventually equal to  $(\mathcal{T}^{(r)})^\perp$ . Focus on one of the  $(\mathcal{T}^{(r)})^\perp$ , for any unit vector (in terms of the  $\|\cdot\|_{\mathcal{A}_r}$  norm)  $\mathbf{a} \in (\mathcal{T}^{(r)})^\perp$ , we have

$$|\langle \mathbf{a}, \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*) \rangle| < \lambda_r - \epsilon. \quad (24)$$

Since the sequence generated by our algorithm converges to the global optimum (as proved in Theorem 1), there exists a  $T$  such that

$$\|\nabla\mathcal{L}(\bar{\boldsymbol{\theta}}_t) - \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*)\|_{\mathcal{A}_r}^* < \epsilon,$$

combining with (24) we have

$$\begin{aligned} |\langle \mathbf{a}, \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}_t) \rangle| &\leq |\langle \mathbf{a}, \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*) \rangle| + |\langle \mathbf{a}, \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}_t) - \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*) \rangle| \\ &< \lambda_r - \epsilon + \|\nabla\mathcal{L}(\bar{\boldsymbol{\theta}}_t) - \nabla\mathcal{L}(\bar{\boldsymbol{\theta}}^*)\|_{\mathcal{A}_r}^* \\ &< \lambda_r \end{aligned} \quad (25)$$

for all  $t > T$ . Now we consider two cases:

1. If  $\langle \mathbf{a}, \boldsymbol{\theta}_{t-1} \rangle \neq 0$ , then  $\mathbf{a} \notin \mathcal{S}_{\text{free}}^{(r)}$  at the  $t$ -th iteration. Since we assume subproblems are exactly solved, and  $\mathbf{a} \in \mathcal{S}_{\text{free}}^{(r)}$ , by the optimality condition  $|\langle \nabla\mathcal{L}(\boldsymbol{\theta}_t), \mathbf{a} \rangle| < \lambda$  implies that  $\langle \boldsymbol{\theta}_t, \mathbf{a} \rangle = 0$ .
2. If  $\langle \mathbf{a}, \boldsymbol{\theta}_{t-1} \rangle = 0$ , combined with (25) we have  $\langle \boldsymbol{\theta}_t, \mathbf{a} \rangle = 0$ .

Therefore, for all  $\mathbf{a} \in (\mathcal{T}^{(r)})^\perp$ , we have  $\langle \boldsymbol{\theta}_t, \mathbf{a} \rangle = 0$ , which implies  $(\mathcal{T}^{(r)})^\perp \subset \mathcal{S}_{\text{fixed}}$ . On the other hand, by definition  $(\mathcal{T}^{(r)}) \cap \mathcal{S}_{\text{fixed}} = \phi$ , thus  $\mathcal{T}^{(r)} = \mathcal{S}_{\text{free}}$  and  $(\mathcal{T}^{(r)})^\perp = \mathcal{S}_{\text{fixed}}$ .  $\square$

### 7.4 Proof of Theorem 3

*Proof.* As shown in Theorem 2, there exists a  $T$  such that  $\mathcal{S}_{\text{free}} = \text{span}(\{\mathbf{a} \mid \langle \mathbf{a}, \boldsymbol{\Theta}^* \rangle\})$  after  $t > T$ .

Next we show that after finite iterations the line search step size will be 1. For simplicity, let  $\bar{\mathcal{L}}(\boldsymbol{\theta}) = \mathcal{L}(\sum_{r=1}^k \boldsymbol{\theta}^{(r)})$  and  $\bar{\mathcal{R}}(\boldsymbol{\theta}) = \sum_{r=1}^k \|\boldsymbol{\theta}_r\|_{\mathcal{A}_r}$ , and  $F(\boldsymbol{\theta}) = \bar{\mathcal{L}}(\boldsymbol{\theta}) + \bar{\mathcal{R}}(\boldsymbol{\theta})$ . Since  $\nabla^2\mathcal{L}(\bar{\boldsymbol{\theta}})$  is Lipschitz continuous, we have

$$\bar{\mathcal{L}}(\boldsymbol{\theta} + \mathbf{d}) \leq \bar{\mathcal{L}}(\boldsymbol{\theta}) + \langle \nabla\bar{\mathcal{L}}(\boldsymbol{\theta}), \mathbf{d} \rangle + \frac{1}{2}\mathbf{d}^T \nabla^2\bar{\mathcal{L}}(\boldsymbol{\theta})\mathbf{d} + \frac{1}{6}\eta\|\mathbf{d}\|^3.$$

Plug in into the objective function we have

$$\begin{aligned} F(\boldsymbol{\theta} + \mathbf{d}) &\leq \bar{\mathcal{L}}(\boldsymbol{\theta}) + \bar{\mathcal{R}}(\boldsymbol{\theta}) + \langle \nabla\bar{\mathcal{L}}(\boldsymbol{\theta}), \mathbf{d} \rangle + \\ &\quad \frac{1}{2}\mathbf{d}^T \nabla^2\bar{\mathcal{L}}(\boldsymbol{\theta})\mathbf{d} + \frac{1}{6}\eta\|\mathbf{d}\|^3 + \bar{\mathcal{R}}(\boldsymbol{\theta} + \mathbf{d}) - \bar{\mathcal{R}}(\boldsymbol{\theta}) \\ &\leq F(\boldsymbol{\theta}) + \delta + \frac{1}{2}\mathbf{d}^T \nabla^2\bar{\mathcal{L}}(\boldsymbol{\theta})\mathbf{d} + \frac{1}{6}\eta\|\mathbf{d}\|^3 \end{aligned} \quad (26)$$

To further bound (26), we first show the following Lemma:

**Lemma 3.** Let  $\mathbf{d}^*$  be the optimal solution of (20), then

$$\delta = \langle \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}), \mathbf{d}^* \rangle + \bar{\mathcal{R}}(\boldsymbol{\theta} + \mathbf{d}^*) - \bar{\mathcal{R}}(\boldsymbol{\theta}) \leq -(\mathbf{d}^*)^T \nabla^2 \bar{\mathcal{L}}(\boldsymbol{\theta}) \mathbf{d}^*.$$

*Proof.* Since  $\mathbf{d}^*$  is the optimal solution of (20), for any  $\alpha < 1$  we have

$$\begin{aligned} \langle \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}) + \mathbf{d}^* \rangle + \frac{1}{2}(\mathbf{d}^*)^T H \mathbf{d}^* + \bar{\mathcal{R}}(\boldsymbol{\theta} + \mathbf{d}^*) &\leq \langle \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}), \alpha \mathbf{d}^* \rangle + \frac{1}{2} \alpha^2 (\mathbf{d}^*)^T H \mathbf{d}^* + \bar{\mathcal{R}}(\boldsymbol{\theta} + \alpha \mathbf{d}^*) \\ &\leq \alpha \langle \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}), \mathbf{d}^* \rangle + \frac{1}{2} \alpha^2 (\mathbf{d}^*)^T H \mathbf{d}^* + \alpha \bar{\mathcal{R}}(\boldsymbol{\theta} + \mathbf{d}^*) + (1 - \alpha) \bar{\mathcal{R}}(\boldsymbol{\theta}), \end{aligned}$$

where we use  $H$  to denote the exact Hessian  $\nabla^2 \bar{\mathcal{L}}(\boldsymbol{\theta})$  and the second inequality is by the convexity of  $\bar{\mathcal{R}}$  since we consider all the atomic norm  $\|\cdot\|_{\mathcal{A}}$  to be convex. Therefore

$$(1 - \alpha) (\langle \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}), \mathbf{d}^* \rangle + \bar{\mathcal{R}}(\boldsymbol{\theta} + \mathbf{d}^*) - \bar{\mathcal{R}}(\boldsymbol{\theta})) + \frac{1}{2} (1 - \alpha^2) (\mathbf{d}^*)^T H \mathbf{d}^* \leq 0.$$

Dividing both sides by  $(1 - \alpha)$  we get

$$\delta + \frac{1}{2} (1 - \alpha) (\mathbf{d}^*)^T H \mathbf{d}^* \leq 0,$$

therefore  $\delta \leq -(\mathbf{d}^*)^T H \mathbf{d}^*$ . □

Combining Lemma 3 and (26) we get

$$F(\boldsymbol{\theta} + \mathbf{d}^*) \leq F(\boldsymbol{\theta}) + \frac{\delta}{2} + \frac{1}{6} \eta \|\mathbf{d}^*\|^3.$$

Furthermore, considering  $\boldsymbol{\theta}$  in the level set, we can define  $M$  to be the largest eigenvalue of Hessians, and thus

$$F(\boldsymbol{\theta} + \mathbf{d}^*) \leq F(\boldsymbol{\theta}) + \left(\frac{1}{2} - \frac{1}{6} \eta M^2 \|\mathbf{d}^*\|\right) \delta.$$

Since  $\mathbf{d} \rightarrow 0$  as  $\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*$ , we can find an  $\epsilon$ -ball around  $x^*$  such that  $\frac{1}{2} - \frac{1}{6} \eta M^2 \|\mathbf{d}^*\| > \sigma$ , and  $\delta < 0$ , thus line search will be satisfied with step size equals to 1.

Based on the above proofs, when  $\boldsymbol{\theta}$  is closed enough to  $\boldsymbol{\theta}^*$ ,  $\mathcal{S}_{\text{free}} = \text{span}(\{\mathbf{a} \mid \langle \mathbf{a}, \boldsymbol{\theta}^* \rangle\})$  and the step size  $\alpha = 1$ . Finally we have to explore the structure of dirty model. Since we have  $k$  parameter sets  $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(k)}$ , the Hessian has a block structure presented in (7). Even when  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$  is strongly convex, the rank of the Hessian matrix is at most  $O(n)$ , while there are totally  $nk$  variables or rows. However, our main observation is that when  $\nabla^2 \mathcal{L}(\boldsymbol{\theta})$  is positive definite, the Hessian  $H \in \mathcal{R}^{pk \times pk}$  has a fixed null space:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}([I, I, \dots, I][\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(k)}]^T) = \mathcal{L}(E\boldsymbol{\theta}),$$

where  $E = [I, I, \dots, I]$ . The null space of  $H$  is always the null space of  $E$  when  $\mathcal{L}$  is strongly convex itself. The following theorem (Theorem 2 in [28]) can then be applied:

**Lemma 4.** Let  $F(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + h(\boldsymbol{\theta})$  where  $g(\boldsymbol{\theta})$  has a constant null space  $\mathcal{T}^\perp$  and is strongly convex in the subspace  $\mathcal{T}$ , and has Lipschitz continuous second order derivative  $\nabla^2 g(\boldsymbol{\theta})$ . If we apply a proximal Newton method (BCGD-block1 with exact Hessian and step size 1) to minimize  $F(\boldsymbol{\theta})$ , then

$$\|\mathbf{z}_{t+1} - \mathbf{z}^*\| \leq \frac{L_H}{2m} \|\mathbf{z}_t - \mathbf{z}^*\|^2,$$

where  $\mathbf{z}^* = \text{proj}_{\mathcal{T}}(\boldsymbol{\theta}^*)$ ,  $\mathbf{z}_t = \text{proj}_{\mathcal{T}}(\mathbf{z}_t)$ , and  $L_H$  is the Lipschitz constant for  $\nabla^2 g(\boldsymbol{\theta})$ .

In our case,  $\text{proj}_{\mathcal{T}}([\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(k)}]^T) = \sum_{i=1}^k \boldsymbol{\theta}^{(i)}$ , therefore we have

$$\left\| \sum_{i=1}^k \boldsymbol{\theta}_{t+1}^{(i)} - \boldsymbol{\theta}^* \right\| \leq C \left\| \sum_{i=1}^k \boldsymbol{\theta}_t^{(i)} - \boldsymbol{\theta}^* \right\|^2,$$

therefore  $\bar{\boldsymbol{\theta}} = \sum_{i=1}^k \boldsymbol{\theta}^{(i)}$  has an asymptotic convergence rate. □

## 7.5 Dual of latent GMRF

The problem (4) can be rewritten by

$$\min_{S, L: L \succeq 0, S-L \succ 0} -\log \det(S-L) + \langle S-L, \Sigma \rangle + \max_{Z: \|Z\|_\infty \leq \alpha} \text{trace}(ZS) + \max_{P: \|P\|_2 \leq \beta} \text{trace}(LP),$$

where  $\|Z\|_\infty = \max_{i,j} |Z_{ij}|$  and  $\|P\|_2 = \sigma_1(P)$  is the induced two norm. We then interchange min and max to get

$$\begin{aligned} \min_{L, S: S-L \succ 0, L \succeq 0} \max_{Z, P: \|Z\|_\infty \leq \alpha, \|P\|_2 \leq \beta} & -\log \det(S-L) + \langle S-L, \Sigma \rangle + \text{trace}(ZS) + \text{trace}(LP) \\ & \equiv g(Z, P, L, S). \end{aligned} \quad (27)$$

Assume we do not have the constraint  $L \succeq 0$ , then the minimizer will satisfy

$$\begin{aligned} \nabla_L g(Z, P, L, S) &= -(S-L)^{-1} + \Sigma + Z = 0 \\ \nabla_S g(Z, P, L, S) &= -(S-L)^{-1} - \Sigma + P = 0. \end{aligned}$$

Therefore we have

$$Z = -P \text{ and } S-L = (\Sigma + Z)^{-1}. \quad (28)$$

Combining (28) and (27) we get the dual problem

$$\min_{\Sigma+Z \geq 0} \log \det(\Sigma + Z) + p \text{ s.t. } \|Z\|_\infty \leq \alpha, \|Z\|_2 \leq \beta.$$

## 7.6 Alternating Minimization approach for latent GMRF

Another way to solve the latent GMRF problem is to directly applying an alternating minimization scheme to solve (4). The algorithm iteratively fix one of the  $S, L$  and update the other. We can still conduct the same active subspace selection technique mentioned in Section 3 in this algorithm. However, this alternating minimization approach can achieve at most linear convergence rate, while our algorithm can achieve super-linear convergence. In this section, we show the detail implementation for this algorithm — ALM(active), and show the comparison results with the Quic & Dirty algorithm (Algorithm 1). The results in Figure 2 shows that our proximal Newton method is faster in terms of the convergence rate.

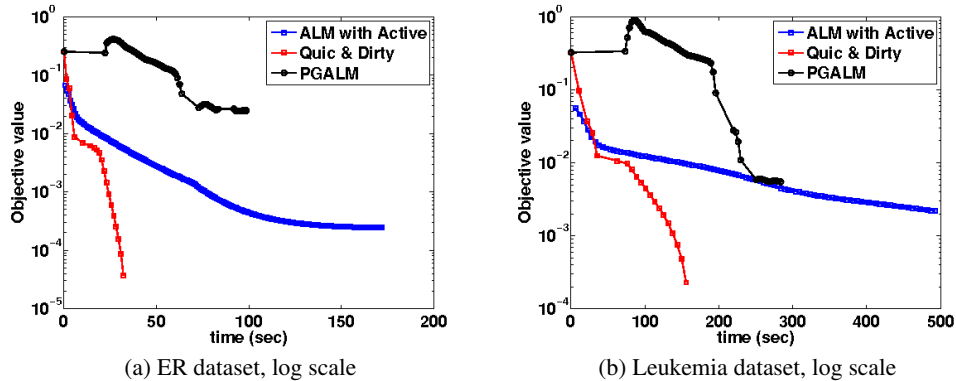


Figure 2: Comparison between QUIC and Alternating Minimization (AM) on gene expression datasets. Note that we implement the active subspace selection approach on both algorithm, so two algorithms have similar speed per iteration. However, we observe that QUIC is much more efficient in terms of final convergence rate.

## 7.7 Details on the implementation of our multi-task solver

We consider the multi-task learning problem. Assume we have  $k$  tasks, each with samples  $X^{(r)} \in \mathcal{R}^{d, n_r}$  and labels  $\mathbf{y}^{(r)} \in \mathcal{R}^{n_r}$ . The goal of multi-task learning is to estimate the model  $W \in \mathcal{R}^{d \times k}$ ,

where each column of  $W$ , denoted by  $\mathbf{w}^{(r)}$ , is the model for the  $r$ -th task. A dirty model has been proposed in [15] to estimate  $W$  by  $S + B$ , where

$$(S, B) = \operatorname{argmin}_{S, B \in \mathcal{R}^{d \times k}} \|\mathbf{y}^{(k)} - X^{(k)}(\mathbf{s}^{(k)} + \mathbf{b}^{(k)})\|^2 + \lambda_S \|S\|_1 + \lambda_B \|B\|_{1,2}, \quad (29)$$

where  $\mathbf{s}^{(k)}, \mathbf{b}^{(k)}$  are the  $k$ -th column of  $S, B$ . It was shown in [15] that the combination of sparse and group sparse regularization yields better performance both in theory and in practice.

Instead of considering the squared-loss problem in (29), we further consider optimization problem minimizing the logistic loss:

$$(S, B) = \operatorname{argmin}_{S, B \in \mathcal{R}^{d \times k}} \sum_{r=1}^k \left( \sum_{i=1}^{n_r} \ell_{\text{logistic}}(\mathbf{y}_i^{(k)}, (\mathbf{s}^{(k)} + \mathbf{b}^{(k)})^T \mathbf{x}_i^{(k)}) \right) + \lambda_S \|S\|_1 + \lambda_B \|B\|_{1,2}, \quad (30)$$

where  $\ell_{\text{logistic}}(y, a) = \log(1 + e^{-ya})$ . This loss function is more suitable for the classification case, as shown in the later experiments.

Let  $W = S + B$ , and we define  $\mathbf{w}^{(r)}$  to be the  $r$ -th column of  $W$ , then the Hessian and gradient of  $\mathcal{L}(\cdot)$  can be computed by

$$\nabla_{\mathbf{w}^{(r)}} \sum_{i=1}^{n_r} (\sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle) - 1) y_i^{(r)} \mathbf{x}_i^{(r)}, \quad \nabla_{\mathbf{w}^{(r)}, \mathbf{w}^{(r)}}^2 \mathcal{L}(W) = (X^{(r)})^T D^{(r)} X^{(r)},$$

where  $\sigma(a) = 1/(1 + e^{-a})$  and  $D^{(r)}$  is a diagonal matrix with  $D_{ii}^{(r)} = \sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle)$  for all  $i = 1, \dots, n_r$ . Note that the Hessian of  $\mathcal{L}(W)$  is a block-diagonal matrix, so  $\nabla_{\mathbf{w}^{(r)}, \mathbf{w}^{(t)}}^2 \mathcal{L}(W) = 0$ . Note also that to form the quadratic approximation at  $W$ , we only need to compute  $\sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle)$  for all  $i, r$ .

For the sparse-structured parameter component, we select a subset of variables in  $S$  to update as in the previous example. For the group-sparse structured component  $B$ , we select a subset of ‘‘rows’’ in  $B$  to update, following (15). To solve the quadratic approximation subproblem, we again use coordinate descent to minimize with respect to the sparse  $S$  component. For the group-sparse component  $B$ , we use block coordinate descent, where each time we update variables in one group (one row) using the trust region approach described in [20]. Since  $\mathcal{S}_{\text{free}}$  contains only a small subset of blocks, the block coordinate descent can focus on this subset and becomes very efficient.

**Algorithm.** We first derive the quadratic approximation for the logistic loss function (30). Let  $W = S + B$ , the gradient for the loss function  $\mathcal{L}(W)$  can be written as

$$\nabla_{\mathbf{w}^{(r)}} \sum_{i=1}^{n_r} (\sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle) - 1) y_i^{(r)} \mathbf{x}_i^{(r)},$$

where  $\sigma(a) = 1/(1 + e^{-a})$ . The Hessian of  $\mathcal{L}(W)$  is a block-diagonal matrix, where each  $d \times d$  block corresponds to variables in one task, i.e.,  $\mathbf{w}^{(r)}$ . Let  $H \in \mathcal{R}^{kd \times kd}$  Hessian matrix, each  $d \times d$  block can be written as

$$\nabla_{\mathbf{w}^{(r)}, \mathbf{w}^{(r)}}^2 \mathcal{L}(W) = (X^{(r)})^T D^{(r)} X^{(r)},$$

where  $D^{(r)}$  is a diagonal matrix with  $D_{ii}^{(r)} = \sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle)$  for all  $i = 1, \dots, n_r$ .

Therefore, to form the quadratic approximation of the current solution, the only computation required is to compute  $\sigma(y_i^{(r)} \langle \mathbf{w}^{(r)}, \mathbf{x}_i^{(r)} \rangle)$  for all  $i, r$ .

For the lasso part, we select a subset of variables in  $S$  to update, according to the subspace selection criterion described 3.2. For the group lasso, we select a subset of ‘‘rows’’ in  $B$  to update, as described in 3.2.

For the lasso part, we apply a coordinate descent solver for solving the subproblem. Notice that for the Lasso part each column of  $S$  forms a subproblems:

$$\min_{\delta \in \mathbb{R}^d} \frac{1}{2} \delta^T H^{(r)} \delta + \mathbf{g}^{(r)} \delta + \|\mathbf{s}^{(r)} + \delta\|,$$

where  $H^{(r)} = (X^{(r)})^T D^{(r)} X^{(r)}$  and  $\mathbf{g}^{(r)} = \nabla_{\mathbf{s}^{(r)}} \mathcal{L}(W)$ . The  $k$  subproblems are independent to each other, so we can solve them independently. For each subproblem, we apply the coordinate

descent approach described in [29] to solve it. When update the coordinate  $\delta_i$ , the key computation is to compute the current gradient  $H^{(r)}\delta + \mathbf{g}^{(r)}$ . Directly computing this is expensive, however, we can maintain  $\mathbf{p} = D^{(r)}X^{(r)}\delta$  during the updates, and then compute  $H^{(r)}\delta = (X_{:,i}^{(r)})^T \mathbf{p}$ , which only takes  $O(\|X_{:,i}^{(r)}\|_0)$  flops.

For solving the group lasso problem, we cannot solve each column independently because the regularization is grouping each row of  $B$ . We apply a block coordinate descent method, where each time only one row of  $B$  is updated. Let  $\delta \in \mathbb{R}^k$  denote the update on the  $i$ -th row of  $B$ , the subproblem with respect to  $\delta$  can be written as

$$\frac{1}{2} \sum_{r=1}^k \gamma_r (\delta_j)^2 + \mathbf{g}^T \delta + \lambda \|\delta + \bar{\mathbf{w}}\|, \quad (31)$$

where  $\bar{\mathbf{w}}$  is the  $i$ -th row of  $W$ ;  $\gamma_r = H_{ii}^{(r)}$  and  $\mathbf{g}_r = \nabla_{W_{ir}} \mathcal{L}(W)$  can be precomputed and will not change during the update.

By taking the gradient of the subproblem (31), we can see that

$$\delta = -\bar{\mathbf{w}} + \begin{cases} 0 & \text{if } \|\mathbf{g} - \sum_{r=1}^k \gamma_r \bar{\mathbf{w}}_r^2\| \leq \lambda \\ -(\Gamma + \frac{\lambda}{\|\bar{\mathbf{w}} + \delta\|} I)^{-1} \mathbf{g} & \text{if } \|\mathbf{g} - \sum_{r=1}^k \gamma_r \bar{\mathbf{w}}_r^2\| > \lambda. \end{cases} \quad (32)$$

For the second case, the closed form solution exists when  $\Gamma = I$ . However, this is not true in general. Instead, we use the iterative trust-region solver proposed in [20] to solve the subproblem, where each iteration of the Newton root finding algorithm only takes  $O(k)$  time. Therefore, the computational bottleneck is to compute  $\mathbf{g}$  in (31). In our case, similar to the Lasso subproblem, we can maintain  $\mathbf{p}^{(r)} = D^{(r)}X^{(r)}\delta^{(r)}$  for each  $r = 1, \dots, k$  in memory, where  $\delta^{(r)}$  is the  $r$ -th column of change in  $W$ . The gradient can then be computed by  $\mathbf{g} = H^{(r)}\delta^{(r)} = X^{(r)}\delta^{(r)}$ , therefore the time complexity is  $O(\bar{n})$  for each coordinate update, where  $\bar{n}$  is number of nonzero for each column of  $X^{(r)}$ .

## 7.8 Proof of Proposition 1

To prove (a), first we expand the sub-differential

$$\langle \mathbf{a}, \partial_{\boldsymbol{\theta}^{(r)}} F(\boldsymbol{\theta}) \rangle = \langle \mathbf{a}, \partial_{\boldsymbol{\theta}^{(r)}} \mathcal{L}(\bar{\boldsymbol{\theta}}) + \lambda_r \partial_{\boldsymbol{\theta}^{(r)}} \|\boldsymbol{\theta}^{(r)}\|_{\mathcal{A}_r} \rangle = \langle \mathbf{a}, G \rangle + \lambda_r \langle \mathbf{a}, \rho \rangle \quad \text{for } \rho \in \partial_{\boldsymbol{\theta}^{(r)}} \|\boldsymbol{\theta}^{(r)}\|_{\mathcal{A}_r}$$

and now using the properties of decomposable norms, we calculate

$$\begin{aligned} |\langle \mathbf{a}, \rho \rangle| &= |\langle \mathbf{a}, \Pi_{\mathcal{T}_r^\perp} \rho \rangle| \\ &\leq \|\mathbf{a}\|_{\mathcal{A}_r} \|\Pi_{\mathcal{T}_r^\perp} \rho\|_{\mathcal{A}_r}^* \\ &\leq 1 \end{aligned}$$

hence  $\langle \mathbf{a}, G \rangle + \lambda_r \langle \mathbf{a}, \rho \rangle \in \langle \mathbf{a}, G \rangle - \lambda_r, [\langle \mathbf{a}, G \rangle + \lambda_r]$  and the result is shown. In fact, it is not hard to see that every element of the set can be written as  $\langle \mathbf{a}, \rho \rangle$  for some  $\rho \in \partial_{\boldsymbol{\theta}^{(r)}} \|\boldsymbol{\theta}^{(r)}\|_{\mathcal{A}_r}$ .

To prove (b), note that the optimality condition on  $\sigma$  is that  $\sigma^*$  will satisfy  $0 \in \partial_\sigma F(\boldsymbol{\theta} + \sigma^* \mathbf{a})$  and by the chain rule,  $\partial_\sigma F(\boldsymbol{\theta} + \sigma \mathbf{a}) = \langle \mathbf{a}, \partial_{\boldsymbol{\theta}^{(r)}} F(\boldsymbol{\theta} + \sigma \mathbf{a}) \rangle$ . If  $\langle \mathbf{a}, G \rangle \leq \lambda_r$ , then by part (a) we see that 0 is in the sub-differential of  $\partial_\sigma F(\boldsymbol{\theta})$  and hence  $\sigma = 0$  is an optimal point. If  $F$  is strongly convex,  $\sigma = 0$  is the unique optimal point.

## 7.9 Proof of Proposition 2

By the definition of proximal operator, in the optimal solution  $-G \in \partial_{\boldsymbol{\theta}^{(r)}} \|\text{prox}_{\lambda_r}(G)\|_{\mathcal{A}_r}$ . For any  $\mathbf{a} \in \mathcal{S}_{fixed}^{(r)}$ ,  $\mathbf{a} \in \mathcal{T}(\text{prox}_{\lambda_r}^{(r)}(G))^\perp$ , thus  $|\langle G, \mathbf{a} \rangle| < \lambda_r$ . Next we consider the projection of gradient to  $\mathcal{S}_{fixed}^{(r)}$ : let  $\rho = \Pi_{\mathcal{S}_{fixed}^{(r)}}(G)$ , then by the previous statement we know  $\|\rho\|_* \leq \lambda_r$ . Then since  $\rho \in \mathcal{T}(\boldsymbol{\theta}^{(r)})^\perp$ , we have  $\rho \in \lambda_r \partial \|\boldsymbol{\theta}^{(r)}\|_{\mathcal{A}_r}$ , therefore constrained to the subspace  $\mathcal{S}_{fixed}^{(r)}$ , 0 belongs to the sub-gradient, which proves Proposition 2.

## 7.10 Active Subspace Selection for Group Lasso Regularization

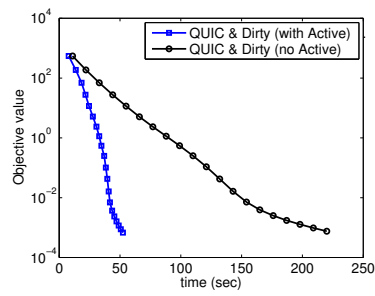


Figure 3: Comparing with/without active subspace selection technique on the RCV1 dataset on the multi-task learning problem with group-lasso regularization and logistic loss. We choose  $\lambda = 10^{-3}$  and the final solution only has 1678 nonzero rows, while there are 22283 rows in total.