

1 First of all, we would thank all the reviewers for their  
 2 careful reading and thoughtful comments. We address  
 3 their comments and questions below. We will prepare our  
 4 final version based on their comments.

5 **Response to Reviewer 1:**

6 - **This method seems to only be applicable to disjoint**  
 7 **groups of coefficients. Many applications require**  
 8 **overlapping groups (e.g.,...).** Does this method trivially  
 9 extend? It does not seem to, at least to me.

10 Answer: Our method is relatively easy to be extended for  
 11 overlapping groups by using overlap norm, which associ-  
 12 ates to the parameter vector a specific decomposition  
 13 [a]. Since we can still use BCD for the overlap norm with  
 14 a slight modification, our method is also available to the  
 15 BCD with the overlapping group.

16 - **This method focuses only on increasing the efficiency**  
 17 **of the block-coordinate descent. Other methods, such**  
 18 **as iterative thresholding for group-sparse inference,**  
 19 **which does not use a full internal block iteration and**  
 20 **instead only performs one step of thresholding per**  
 21 **block, can be used instead (e.g.,...).** Do any of the ben-  
 22 efits transfer? Depending on the number of iterations it takes to compute, these methods could be comparable,  
 23 or be much slower. It would be useful for the authors to compare to these other methods

24 Answer: Yes, our ideas of the upper bound and the candidate group set is applicable to other methods such as the  
 25 iterative thresholding with a slight modification of Eq. (6) and Eq. (8). This is because the computations of the bound  
 26 and the set are based on the difference between the reference parameter vector and the current parameter vector, which  
 27 can also be handled in the iterative thresholding.

28 - **I would like to see comparisons to other**

29 We compared our method with an iterative thresholding-style algorithm. We shows the result for  $\alpha = 0.2$  in Figure 1 as  
 30 a example. Our method outperformed Iterative Thresholding. The similar results have been observed for other  $\alpha$ . We  
 31 will add all the results to the final version of our paper. We thank the reviewer for this constructive comment.

32 **Response to Reviewer 2:**

33 - **a figure giving the ratio of activated variable could be useful to compare the selection process (see fig. 4 in [1])**  
 34 Figure 2 shows the result of the selection process on the boston dataset by following Figure 4 in [1].  $K$  represents the  
 35 number of the main loop of BCD using the upper bound. Our method effectively skips the variables even if  $\lambda$  is small.  
 36 We will add the results of other datasets in the final version of our paper. We appreciate your insightful advice.

37 - **clarify at the beginning that the method only consider non-overlapping groups**

38 We will revise the paper to reflect this comment. We thank the reviewer for your important advice.

39 - **what are some possible extension of the proposed algorithm ?**

40 Our method can be flexibly extended to various types of group structured data. For examples, our method can be  
 41 naturally extended to overlapping groups by using overlap norm [a]. In addition, our method is also applicable for  
 42 tree-structured data since Tree-structured Group Lasso [b] uses a similar equation to Eq. (4) of SGL.

43 **Response to Reviewer 3:**

44 - **Provide a more clear proof of Theorem 2. (Related comment: This proof seems to ignore a major difference in**  
 45 **the group updating order...)**

46 Theorem 2 holds regardless of the group updating order because it guarantees *the converged value of the objective*  
 47 *function* rather than *the converged parameter vector*. Since the problem of SGL is either convex or strongly convex  
 48 depending on the condition of  $X$ , the *optimal value of the objective function* is unique (while *the optimal parameter*  
 49 *vector* may be non-unique). In addition, because we assume that the original BCD converges to the solution in the  
 50 theorem, the BCD using the upper bound also converges even if it starts with the initial parameter resulting from the  
 51 BCD on the candidate group set. We will add this explanation to the paper in detail.

52 - **Derive a convergence rate and compare that with the original BCD algorithm.**

53 Our method performs the first BCD on the candidate group set and then performs the second BCD on all the groups  
 54 using the upper bound. We would like to analyze the first BCD in the future work because our own criterion of Eq. (8)  
 55 is nontrivial for the convergence analysis. However, since the updating order of the second BCD is the same as the  
 56 original BCD, the second BCD achieves at least the same convergence rate as that of the original BCD with the initial  
 57 parameter vector resulting from the first BCD.

58 - **State Lemma 4 in a more clear way. (Related question: Does Lemma 4 mean that the candidate group set**  
 59 **contains \*ALL\* groups whose parameters must be nonzeros?)**

60 No, the candidate group set contains a subset of the nonzero groups *identified by using the lower bound of Lemma 3.2*  
 61 as shown in the proof of Lemma 4. The lower bound confidently identifies groups with *nonzero* vectors while the upper  
 62 bound of Lemma 1 identifies groups with *zero* vectors. We are sorry for the confusion. We will modify Lemma 4 to add  
 63 the aforementioned explanation.

64 [a] L. Jacob, G. Obozinski, J. Vert. Group lasso with overlap and graph lasso. In ICML, 2009.

65 [b] R. Jenatton, J. Mairal, G. Obozinski, F. Bach. Proximal Methods for Sparse Hierarchical Dictionary Learning. In ICML, 2010.

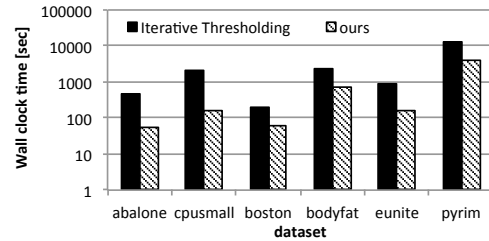


Figure 1: Comparison with iterative thresholding. ( $\alpha = 0.2$ )

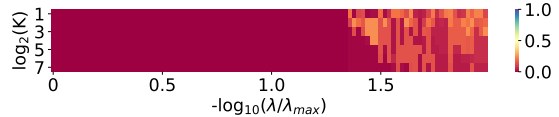


Figure 2: The ratio of active variables during the optimization of our method.  $K$  represents the number of the main loop of BCD using the upper bound. ( $\alpha = 0.2$ )