# Supplementary Materials of "Limited-memory Common-directions Method for Distributed L1-regularized Linear Classification"

Wei-Lin Chiang[*]    Yu-Sheng Li[†]    Ching-pei Lee[‡]    Chih-Jen Lin[§]

## I  Introduction

In this document, we present additional details and more experimental results.

## II  Derivation of the Direction Used in LBFGS

By using only the information from the last $m$ iterations, the definition of $B_k$ in BFGS becomes the following in LBFGS.

(II.1)
$$B_k = V_{k-1}^T \cdots V_{k-m}^T B_0^k V_{k-m} \cdots V_{k-1}$$
$$+ \rho_{k-m} V_{k-1}^T \cdots V_{k-m+1}^T \boldsymbol{s}_{k-m} \boldsymbol{s}_{k-m}^T V_{k-m+1} \cdots V_{k-1}$$
$$+ \cdots + \rho_{k-1} \boldsymbol{s}_{k-1} \boldsymbol{s}_{k-1}^T.$$

Note that in BFGS, $B_0^k$ is a fixed matrix, but in LBFGS, $B_0^k$ can change with $k$, provided its eigenvalues are bounded in a positive interval over $k$. A common choice is

(II.2)
$$B_0^k = \frac{\boldsymbol{s}_{k-1}^T \boldsymbol{u}_{k-1}}{\boldsymbol{u}_{k-1}^T \boldsymbol{u}_{k-1}} I.$$

By expanding (II.1), $\boldsymbol{d}_k$ can be efficiently obtained by $\mathcal{O}(m)$ vector operations as shown in Algorithm II. The overall procedure of LBFGS is summarized in Algorithm I.

## III  More Details of Limited-memory Common-directions Method

A sketch of the procedure for L1-regularized problems is in Algorithm III.

## IV  Line Search in Algorithms for L2-regularized Problems

Here we present the trick mentioned in Section 2.3 in the paper. At each line search iteration, we obtain $\boldsymbol{w}^T \boldsymbol{x}_i, \boldsymbol{d}^T \boldsymbol{x}_i, \forall i$ first, and then use $\mathcal{O}(l)$ cost to calculate

$$(\boldsymbol{w} + \alpha \boldsymbol{d})^T \boldsymbol{x}_i = \boldsymbol{w}^T \boldsymbol{x}_i + \alpha \boldsymbol{d}^T \boldsymbol{x}_i, \forall i.$$

---

[*]National Taiwan University. `b02902056@ntu.edu.tw`
[†]National Taiwan University. `b03902086@ntu.edu.tw`
[‡]University of Wisconsin-Madison. `ching-pei@cs.wisc.edu`
[§]National Taiwan University. `cjlin@csie.ntu.edu.tw`

---

**Algorithm I** LBFGS.

1: Given $\boldsymbol{w}_0$, integer $m > 0$, and $\beta, \gamma \in (0, 1)$
2: $\boldsymbol{w} \leftarrow \boldsymbol{w}_0$
3: **for** $k = 0, 1, 2, \ldots$ **do**
4:    Calculate $\nabla f(\boldsymbol{w})$
5:    Calculate the new direction

$$\boldsymbol{d} \equiv - B \nabla f(\boldsymbol{w}) \approx -\nabla^2 f(\boldsymbol{w})^{-1} \nabla f(\boldsymbol{w})$$

   by using information of the previous $m$ iterations (Algorithm II)
6:    Calculate $\nabla f(\boldsymbol{w})^T \boldsymbol{d}$
7:    $\alpha \leftarrow 1, \boldsymbol{w}^{\text{old}} \leftarrow \boldsymbol{w}$
8:    **while** true **do**
9:       $\boldsymbol{w} \leftarrow \boldsymbol{w}^{\text{old}} + \alpha \boldsymbol{d}$
10:       Calculate the objective value $f(\boldsymbol{w})$ in (2.2)
11:       **if** $f(\boldsymbol{w}) - f(\boldsymbol{w}^{\text{old}}) \leq \gamma \nabla f(\boldsymbol{w}^{\text{old}})^T (\boldsymbol{w} - \boldsymbol{w}^{\text{old}})$
12:          **break**
13:       $\alpha \leftarrow \alpha \beta$
14:    Update $P$ with $\boldsymbol{w} - \boldsymbol{w}^{\text{old}}$ and $\nabla f(\boldsymbol{w}) - \nabla f(\boldsymbol{w}^{\text{old}})$

---

Because $\boldsymbol{w}^T \boldsymbol{x}_i$ can be obtained from the previous iteration, $\boldsymbol{d}^T \boldsymbol{x}_i$ is the only $\mathcal{O}(\#\text{nnz})$ operation needed. The line search cost is reduced from

$$\#\text{line-search steps} \times \mathcal{O}(\#\text{nnz})$$

to

$$1 \times \mathcal{O}(\#\text{nnz}) + \#\text{line-search steps} \times \mathcal{O}(l).$$

This trick is not applicable to L1-regularized problems because the new point is no longer $\boldsymbol{w} + \alpha \boldsymbol{d}$.

## V  More on the Distributed Implementation

**V.1  Complexity.** The distributed implementation as mentioned in Section 5 is shown in Algorithm VI. Then we discuss the complexity below.

$$\frac{(2 + \#\text{line-search steps}) \times \mathcal{O}(\#\text{nnz}) + \mathcal{O}(lm^2) + \mathcal{O}(mn)}{K}$$
$$+ \mathcal{O}(m^3)$$

## VI  More Experiments

The data sets used in this section are shown in Table (I).

**Algorithm II** LBFGS Two-loop recursion.

1: $\boldsymbol{q} \leftarrow -\nabla f(\boldsymbol{w})$
2: **for** $i = k-1, k-2, \ldots, k-m$ **do**
3:   $\alpha_i \leftarrow \boldsymbol{s}_i^T \boldsymbol{q} / \boldsymbol{s}_i^T \boldsymbol{u}_i$
4:   $\boldsymbol{q} \leftarrow \boldsymbol{q} - \alpha_i \boldsymbol{u}_i$
5: $\boldsymbol{r} \leftarrow (\boldsymbol{s}_{k-1}^T \boldsymbol{u}_{k-1} / \boldsymbol{u}_{k-1}^T \boldsymbol{u}_{k-1}) \boldsymbol{q}$
6: **for** $i = k-m, k-m+1, \ldots, k-1$ **do**
7:   $\beta_i \leftarrow \boldsymbol{u}_i^T \boldsymbol{r} / \boldsymbol{s}_i^T \boldsymbol{u}_i$
8:   $\boldsymbol{r} \leftarrow \boldsymbol{r} + (\alpha_i - \beta_i) \boldsymbol{s}_i$
9: $\boldsymbol{d} \leftarrow \boldsymbol{r}$

---

**Algorithm III** Limited-memory common-directions method for L1-regularized problems.

1: **while** true **do**
2:   Compute $\nabla^{\mathrm{P}} f(\boldsymbol{w})$ by (2.6)
3:   Solve the sub-problem (3.21)
4:   Let the direction be $\boldsymbol{d} = P\boldsymbol{t}$
5:   **for** $j = 1, \ldots, n$ **do**
6:     Align $d_j$ with $-\nabla_j^{\mathrm{P}} f(\boldsymbol{w})$ by (2.11)
7:   $\alpha \leftarrow 1$, $\boldsymbol{w}^{\mathrm{old}} \leftarrow \boldsymbol{w}$
8:   **while** true **do**
9:     Calculate $\boldsymbol{w}$ from $\boldsymbol{w}^{\mathrm{old}} + \alpha\boldsymbol{d}$ by (2.12)
10:    **if** $f(\boldsymbol{w}) - f(\boldsymbol{w}^{\mathrm{old}}) \leq \gamma\nabla^{\mathrm{P}} f(\boldsymbol{w}^{\mathrm{old}})^T(\boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}})$
11:      **break**
12:    $\alpha \leftarrow \alpha\beta$
13:   Update $P$ and $XP$

**VI.1 More Results by Using Different $C$ Values.** In Figure (I), we present more results with

$$C = \{0.1C_{\mathrm{Best}}, C_{\mathrm{Best}}, 10C_{\mathrm{Best}}\},$$

where $C_{\mathrm{Best}}$ is the value to achieve the highest cross validation accuracy. The results are similar to $C_{\mathrm{Best}}$ presented in Section 6, but we observe that NEWTON converges slowly in larger $C$ cases.

**VI.2 More Results on Distributed Experiments.** In Figure (II), we present more data sets for the distributed experiments. All settings are the same as in Section 6.

Table (I): Data statistics.

| Data set | #instances | #features | #nonzeros | $C_{\mathrm{Best}}$ |
|---|---|---|---|---|
| real-sim | 72,309 | 20,958 | 3,709,083 | 16 |
| rcv1_test | 677,399 | 47,226 | 49,556,258 | 4 |
| news20 | 19,996 | 1,355,191 | 9,097,916 | 1024 |
| yahoojp | 176,203 | 832,026 | 23,506,415 | 2 |
| url | 2,396,130 | 3.231,961 | 277,058,644 | 8 |
| yahookr | 460,554 | 3,052,939 | 156,436,656 | 4 |
| epsilon | 400,000 | 2,000 | 800,000,000 | 0.5 |
| webspam | 350,000 | 16,609,143 | 1,304,697,446 | 64 |
| KDD2010-b | 19,264,097 | 29,890,096 | 566,345,888 | 0.5 |
| criteo | 45,840,617 | 1,000,000 | 1,787,773,969 | 0.5 |
| avazu-site | 25,832,830 | 999,962 | 387,492,144 | 1 |
| kdd2012 | 149,639,105 | 54,686,452 | 1,646,030,155 | 2 |

**Algorithm IV** A distributed implementation of OWLQN.

1: **for** $k = 0, 1, 2, \ldots$ **do**
2:     Compute $\nabla^{\mathrm{P}} f(\boldsymbol{w})$ by (2.6) and

$$\nabla L(\boldsymbol{w}) = C \bigoplus_{r=1}^{K} (X_{J_r,:})^T \begin{bmatrix} \vdots \\ \xi'(y_i \boldsymbol{w}^T \boldsymbol{x}_i) \\ \vdots \end{bmatrix}_{i \in J_r} .$$

                              ▷ $\mathcal{O}(\#\mathrm{nnz}/K)$; $\mathcal{O}(n)$ comm.
3:     Compute the search direction $\boldsymbol{d}_{\bar{J}_r}, r = 1, \ldots, K$
    by Algorithm V        ▷ $\mathcal{O}(nm/K)$; $\mathcal{O}(m)$ comm.
4:     An *allgather* operation to let each node has

$$\boldsymbol{d} = \begin{bmatrix} \boldsymbol{d}_{\bar{J}_1} \\ \vdots \\ \boldsymbol{d}_{\bar{J}_K} \end{bmatrix}$$

                              ▷ $\mathcal{O}(n/K)$ comm.
5:     **for** $j = 1, \ldots, n$ **do**
6:         Align $d_j$ with $-\nabla_j^{\mathrm{P}} f(\boldsymbol{w})$ by (2.11)
7:     $\alpha \leftarrow 1$, $\boldsymbol{w}^{\mathrm{old}} \leftarrow \boldsymbol{w}$
8:     **while** true **do**
9:         Calculate $\boldsymbol{w}$ from $\boldsymbol{w}^{\mathrm{old}} + \alpha \boldsymbol{d}$ by (2.12) and

$$f(\boldsymbol{w}) = \|\boldsymbol{w}\|_1 + C \bigoplus_{r=1}^{K} \sum_{i \in J_r} \xi(y_i \boldsymbol{w}^T \boldsymbol{x}_i)$$

                              ▷ $\mathcal{O}(\#\mathrm{nnz}/K)$; $\mathcal{O}(1)$ comm.
10:         **if** $f(\boldsymbol{w}) - f(\boldsymbol{w}^{\mathrm{old}}) \le \gamma \nabla^{\mathrm{P}} f(\boldsymbol{w}^{\mathrm{old}})^T (\boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}})$
11:             **break**
12:         $\alpha \leftarrow \alpha \beta$
13:     $\boldsymbol{s}_k \leftarrow \boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}}, \quad \boldsymbol{u}_k \leftarrow \nabla f(\boldsymbol{w}) - \nabla f(\boldsymbol{w}^{\mathrm{old}})$
14:     Remove 1st column of $S$ and $U$ if needed and

$$S \leftarrow \begin{bmatrix} S & \boldsymbol{s}_k \end{bmatrix}, \quad U \leftarrow \begin{bmatrix} U & \boldsymbol{u}_k \end{bmatrix}$$

15:     $\rho_k \leftarrow \bigoplus_{r=1}^{K} (\boldsymbol{u}_k)_{\bar{J}_r}^T (\boldsymbol{s}_k)_{\bar{J}_r}$        ▷ $\mathcal{O}(n/K)$; $\mathcal{O}(1)$
    comm.

---

**Algorithm V** Distributed OWLQN Two-loop recursion

1: **if** $k = 0$  **return** $\boldsymbol{d}_{\bar{J}_r} \leftarrow -\nabla_{\bar{J}_r}^{\mathrm{P}} f(\boldsymbol{w})$
2: **for** $r = 1, \ldots, K$ **do in parallel**
3:     $\boldsymbol{q}_{\bar{J}_r} \leftarrow -\nabla_{\bar{J}_r}^{\mathrm{P}} f(\boldsymbol{w})$                ▷ $\mathcal{O}(n/K)$
4: **for** $i = k-1, k-2, \ldots, k-m$ **do**
5:     Calculate $\alpha_i$ by

$$\alpha_i \leftarrow \frac{\bigoplus_{r=1}^{K} (\boldsymbol{s}_i)_{\bar{J}_r}^T \boldsymbol{q}_{\bar{J}_r}}{\rho_i}$$

                              ▷ $\mathcal{O}(n/K)$; $\mathcal{O}(1)$ comm.
6:     **for** $r = 1, \ldots, K$ **do in parallel**
7:         $\boldsymbol{q}_{\bar{J}_r} \leftarrow \boldsymbol{q}_{\bar{J}_r} - \alpha_i (\boldsymbol{u}_i)_{\bar{J}_r}$              ▷ $\mathcal{O}(n/K)$
8: Calculate

$$\boldsymbol{u}_{k-1}^T \boldsymbol{u}_{k-1} \leftarrow \bigoplus_{r=1}^{K} (\boldsymbol{u}_{k-1})_{\bar{J}_r}^T (\boldsymbol{u}_{k-1})_{\bar{J}_r}$$

                              ▷ $\mathcal{O}(n/K)$; $\mathcal{O}(1)$ comm.
9: **for** $r = 1, \ldots, K$ **do in parallel**
10:     $\boldsymbol{r}_{\bar{J}_r} \leftarrow \frac{\rho_{k-1}}{\boldsymbol{u}_{k-1}^T \boldsymbol{u}_{k-1}} \boldsymbol{r}_{\bar{J}_r}$              ▷ $\mathcal{O}(n/K)$
11: **for** $i = k-m, k-m+1, \ldots, k-1$ **do**
12:     Calculate $\beta_i$ by

$$\beta_i \leftarrow \frac{\bigoplus_{r=1}^{K} (\boldsymbol{u}_i)_{\bar{J}_r}^T \boldsymbol{r}_{\bar{J}_r}}{\rho_i}$$

                              ▷ $\mathcal{O}(n/K)$; $\mathcal{O}(1)$ comm.
13:     **for** $r = 1, \ldots, K$ **do in parallel**
14:         $\boldsymbol{r}_{\bar{J}_r} \leftarrow \boldsymbol{r}_{\bar{J}_r} + (\alpha_i - \beta_i)(\boldsymbol{s}_i)_{\bar{J}_r}$        ▷ $\mathcal{O}(n/K)$
    **return** $\boldsymbol{d}_{\bar{J}_r} \leftarrow \boldsymbol{r}_{\bar{J}_r}, r = 1, \ldots, K$

**Algorithm VI** Distributed limited-memory common-directions method.

---

1: **while** true **do**
2:     Compute $\nabla^{\mathrm{P}} f(\boldsymbol{w})$ by (2.6) and

$$\nabla L(\boldsymbol{w}) = C \bigoplus_{r=1}^{K} (X_{J_r,:})^T \begin{bmatrix} \vdots \\ \xi'(y_i \boldsymbol{w}^T \boldsymbol{x}_i) \\ \vdots \end{bmatrix}_{i \in J_r}.$$

                                         $\triangleright \mathcal{O}(\#\mathrm{nnz}/K); \mathcal{O}(n)$ comm.
3:     Calculate
$$X_{J_r,:} \nabla^{\mathrm{P}} f(\boldsymbol{w})$$

                                         $\triangleright \mathcal{O}(\#\mathrm{nnz}/K)$
4:     Remove 1st column of $P$ and $U$ if needed and

$$P_{J_r,:} \leftarrow \begin{bmatrix} P_{J_r,:} & \nabla^{\mathrm{P}}_{J_r} f(\boldsymbol{w}) \end{bmatrix}$$
$$U_{J_r,:} \leftarrow \begin{bmatrix} U_{J_r,:} & X_{J_r,:} \nabla^{\mathrm{P}} f(\boldsymbol{w}) \end{bmatrix}$$

5:     Calculate

$$(XP)^T D_{\boldsymbol{w}}(XP) = \bigoplus_{r=1}^{K} (U_{J_r,:})^T (D_{\boldsymbol{w}})_{J_r,J_r} U_{J_r,:}$$

                                 $\triangleright \mathcal{O}(lm^2/K), \mathcal{O}(m^2)$ comm.

$$-P^T \nabla^{\mathrm{P}} f(\boldsymbol{w}) = -\bigoplus_{r=1}^{K} (P_{\bar{J}_r,:})^T \nabla^{\mathrm{P}}_{\bar{J}_r} f(\boldsymbol{w})$$

                                $\triangleright \mathcal{O}(mn/K); \mathcal{O}(m)$ comm.
6:     Solve

$$\left( (XP)^T D_{\boldsymbol{w}}(XP) \right) \boldsymbol{t} = -P^T \nabla^{\mathrm{P}} f(\boldsymbol{w})$$

                                        $\triangleright \mathcal{O}(m^3)$
7:     Let the direction be

$$\boldsymbol{d} = P\boldsymbol{t} = \begin{bmatrix} P_{\bar{J}_1,:} \boldsymbol{t}, \ldots, P_{\bar{J}_K,:} \boldsymbol{t} \end{bmatrix}^T$$

                                 $\triangleright \mathcal{O}(mn/K); \mathcal{O}(n/K)$ comm.
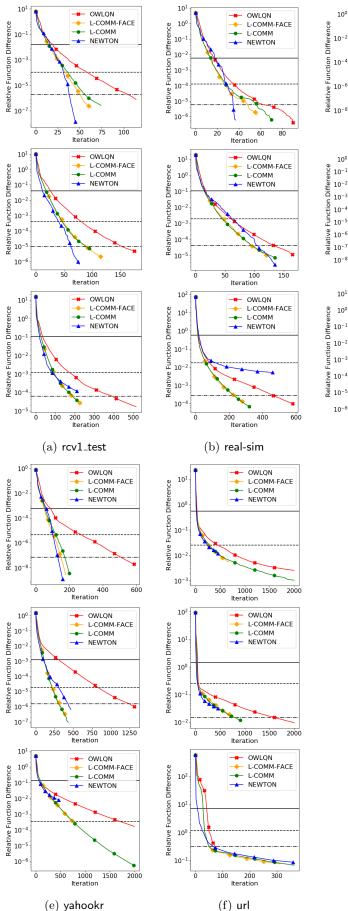8:     **for** $j = 1, \ldots, n$ **do**
9:         Align $d_j$ with $-\nabla^{\mathrm{P}}_j f(\boldsymbol{w})$ by (2.11)
10:     $\alpha \leftarrow 1, \boldsymbol{w}^{\mathrm{old}} \leftarrow \boldsymbol{w}$
11:     **while** true **do**
12:         Calculate $\boldsymbol{w}$ from $\boldsymbol{w}^{\mathrm{old}} + \alpha \boldsymbol{d}$ by (2.12) and
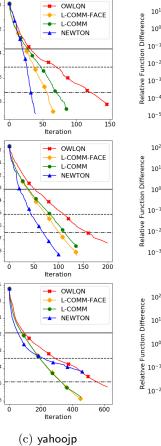
$$f(\boldsymbol{w}) = \|\boldsymbol{w}\|_1 + C \bigoplus_{r=1}^{K} \sum_{i \in J_r} \xi(y_i \boldsymbol{w}^T \boldsymbol{x}_i)$$

                                $\triangleright \mathcal{O}(\#\mathrm{nnz}/K); \mathcal{O}(1)$ comm.
13:         **if** $f(\boldsymbol{w}) - f(\boldsymbol{w}^{\mathrm{old}}) \leq \gamma \nabla^{\mathrm{P}} f(\boldsymbol{w}^{\mathrm{old}})^T (\boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}})$
14:            **break**
15:         $\alpha \leftarrow \alpha\beta$
16:     Remove 1st column of $P$ and $U$ if needed and

$$P \leftarrow \begin{bmatrix} P & \boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}} \end{bmatrix}$$
$$U_{J_r,:} \leftarrow \begin{bmatrix} U_{J_r,:} & X_{J_r,:}(\boldsymbol{w} - \boldsymbol{w}^{\mathrm{old}}) \end{bmatrix}$$

---

(a) rcv1_test     (b) real-sim     (c) yahoojp     (d) news20
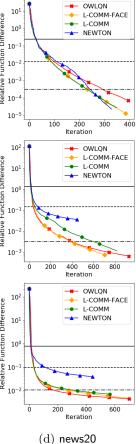
(e) yahookr     (f) url

Figure (I): Comparison of different algorithms with $0.1C_{\text{Best}}$ , $C_{\text{Best}}$ , $10C_{\text{Best}}$ , respectively from top to below for each data set. We show iteration versus the relative difference to the optimal value. Other settings are the same as in Figure 3
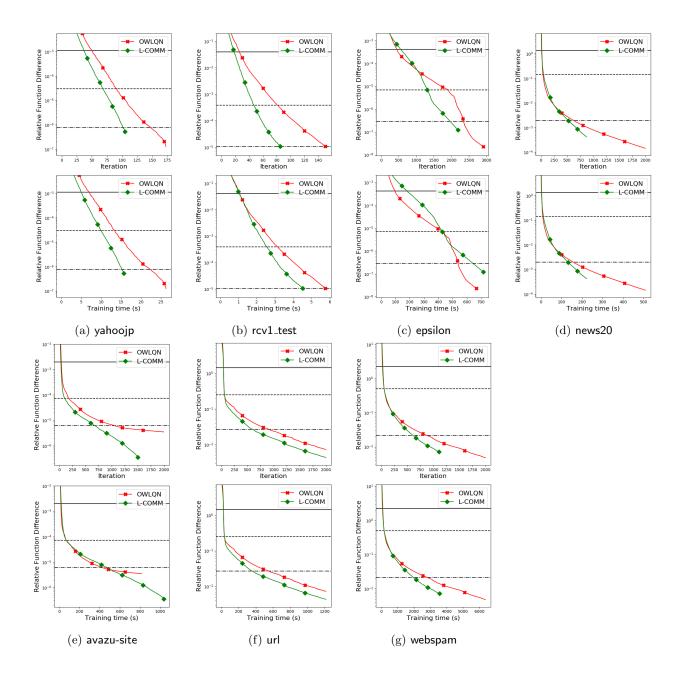
Figure (II): Comparison of different algorithms by using 32 nodes. Upper: iterations. Lower: running time in seconds. Other settings are the same as Figure 4.