

The Common-directions Method for Regularized Empirical Risk Minimization

Po-Wei Wang

Department of Computer Science

National Taiwan University

Taipei 106, Taiwan

B97058@CSIE.NTU.EDU.TW

Ching-pei Lee

Department of Computer Sciences

University of Wisconsin-Madison

Madison, WI 53706-1613, USA

CHING-PEI@CS.WISC.EDU

Chih-Jen Lin

Department of Computer Science

National Taiwan University

Taipei 106, Taiwan

CJLIN@CSIE.NTU.EDU.TW

Editor:

Abstract

State-of-the-art first-order optimization methods are able to achieve optimal linear convergence rate, but cannot obtain quadratic convergence. In this work, we propose a first-order method for regularized empirical risk minimization that exploits the problem structure to efficiently combine multiple update directions to obtain both optimal global linear convergence rate and local quadratic convergence. Experimental results show that our method outperforms state-of-the-art first- and second-order optimization methods in terms of the number of data accesses, while is competitive in training time.

1. Introduction

Consider the general problem

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}), \quad (1)$$

where $f(\mathbf{w})$ is convex and continuously differentiable. Most iterative optimization methods generate iterates by using the first- and second-order derivatives. After a new step is generated, the previous steps and related calculation results are usually considered outdated and thus discarded. In this work, we show that we can reuse previous descent directions in later iterations to reduce computation and to achieve both linear and quadratic convergence.

The idea behind our proposed method is simple. We store all previous directions, and approximately solve the subproblem

$$\min_{\mathbf{t} \in \mathbb{R}^m} f(\mathbf{w} + P\mathbf{t}), \quad (2)$$

where the columns of $P = [\mathbf{p}_1, \dots, \mathbf{p}_m]$ form an orthonormal basis of the span of the previous update directions. The resulting direction $P\mathbf{t}$ is then used to update the current iterate \mathbf{w} for the next iteration. Utilizing all previous directions obtained from the first iteration on seems to be computationally and spatially expensive, especially in the later stage of this method. However, we note that for regularized empirical risk minimization (ERM) of linear models, the method comes with little additional cost by wisely caching the inner products between these \mathbf{p}_i and the training data. We will exemplify how to use this caching strategy to significantly accelerate the computation in our method for ERM. Experimental results show that our method outperforms state of the art.

There are other optimization methods that also reuse previous update directions. For example, the heavy-ball method (Polyak, 1964), the (nonlinear) conjugate gradient method (Fletcher and Reeves, 1964), and Nesterov's accelerated gradient method (Nesterov, 1983, 2013) reuse the previous direction $\mathbf{w}_k - \mathbf{w}_{k-1}$ together with the current gradient to decide the update direction at the k -th iteration. On the other hand, quasi-Newton methods such as BFGS (Dennis and Moré, 1977) and L-BFGS (Liu and Nocedal, 1989) store previous gradients to construct an approximation of the Hessian. More details of these methods and empirical comparisons are provided in later sections.

This paper is organized as follows. We describe the details of our algorithm in Section 2. We prove global linear convergence and local quadratic convergence. Section 3 illustrates important techniques for accelerating the computation when solving ERM problems. We then extend the proposed method to a more general framework and obtain a better linear convergence rate in Section 4. Related works are discussed in Section 5. Experimental results in Section 6 show that our method outperforms state-of-the-art methods empirically in terms of the number of data accesses, and is competitive in terms of training time. We then consider some interesting extensions in Section 7. Section 8 concludes the paper.

Programs used for experiments and a supplementary file including additional results can be found at <http://www.csie.ntu.edu.tw/~cjlin/papers/nheavy/>.

2. Common-directions Method

Before starting the description of our algorithm, we first specify the type of problems we consider. For all our theoretical results except Theorem 4.2, we require that $f(\mathbf{w})$ satisfies Assumption 2.1.

Assumption 2.1 *The objective function $f(\mathbf{w})$ is differentiable, ρ Lipschitz smooth and σ strongly convex with some constants $\rho \geq \sigma > 0$.¹ That is,*

$$\|\nabla f(\mathbf{u}) - \nabla f(\mathbf{v})\| \leq \rho \|\mathbf{u} - \mathbf{v}\|, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n, \quad (3)$$

and

$$f(\mathbf{u}) - f(\mathbf{v}) - \nabla f(\mathbf{v})^\top (\mathbf{u} - \mathbf{v}) \geq \frac{\sigma}{2} \|\mathbf{u} - \mathbf{v}\|^2, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n. \quad (4)$$

For some theoretical results, we further need that $f(\mathbf{w})$ satisfies Assumption 2.2.

1. The requirement $\rho \geq \sigma$ comes from the conditions (3) and (4).

Algorithm 1: A framework for the common-directions method

```

Given  $\mathbf{w}_0$ , compute  $\nabla f(\mathbf{w}_0)$ ;
 $P = [\nabla f(\mathbf{w}_0)/\|\nabla f(\mathbf{w}_0)\|]$ ;
for  $k=0,1,\dots$  do
    if  $\nabla f(\mathbf{w}_k) = \mathbf{0}$ , or some stopping condition is satisfied then
        | stop;
    end
    Let  $\mathbf{t}_k = -(P^\top \nabla^2 f(\mathbf{w}_k) P)^{-1} P^\top \nabla f(\mathbf{w}_k)$ ;
    Backtracking line search on  $f(\mathbf{w}_k + \theta P \mathbf{t}_k)$  by Algorithm 2 to obtain  $\theta_k$ ;
     $\mathbf{w}_{k+1} = \mathbf{w}_k + \theta_k P \mathbf{t}_k$ ;
    Compute  $\nabla f(\mathbf{w}_{k+1})$ ;
    Let  $\mathbf{p} = \nabla f(\mathbf{w}_{k+1}) - P(P^\top \nabla f(\mathbf{w}_{k+1}))$ ;
    if  $\mathbf{p} \neq \mathbf{0}$  then
        |  $P = [P, \mathbf{p}/\|\mathbf{p}\|]$ ;
    end
end

```

Algorithm 2: Backtracking line search

```

Given  $\beta \in (0, 1)$ ,  $\lambda > 0$ , the current iterate  $\mathbf{w}$ , a descent direction  $P \mathbf{t}$ ;
Let  $\theta := 1$ ;
while  $f(\mathbf{w}) - f(\mathbf{w} + \theta P \mathbf{t}) < \frac{\lambda}{2} \theta^2 \|P \mathbf{t}\|^2$  do
    |  $\theta := \beta \theta$ ;
end

```

Assumption 2.2 *The objective function $f(\mathbf{w})$ is twice-differentiable, and the Hessian of $f(\mathbf{w})$ is M Lipschitz continuous with some constant $M > 0$. That is,*

$$\|\nabla^2 f(\mathbf{u}) - \nabla^2 f(\mathbf{v})\| \leq M \|\mathbf{u} - \mathbf{v}\|, \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^n.$$

When minimizing a function, we usually only work on derivatives of the current iterate. For example, a gradient descent step can be derived from searching along the negative gradient direction, and a Newton step can be expressed by the current gradient and Hessian. However, after obtaining the update direction, stale derivatives are usually discarded, which can be inefficient. A simple idea is to store all previous directions \mathbf{p}_j in $P = [\mathbf{p}_1, \dots, \mathbf{p}_m]$ and consider working on the linear span of \mathbf{p}_j . If possible, we search for the best linear combination $P \mathbf{t}$ such that $f(\mathbf{w} + P \mathbf{t})$ is the smallest; see the sub-problem (2). However, in practice (2) may be difficult to be solved, so following past optimization development, we consider the second-order Taylor approximation of $f(\mathbf{w})$ as

$$\tilde{f}(\mathbf{u} \mid \mathbf{w}) \equiv f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{u} - \mathbf{w}) + \frac{1}{2} (\mathbf{u} - \mathbf{w})^\top \nabla^2 f(\mathbf{w}) (\mathbf{u} - \mathbf{w}), \quad (5)$$

where $\nabla^2 f(\tilde{\mathbf{w}})$ is the Hessian at $\tilde{\mathbf{w}}$. Then, a descent direction $\mathbf{d} = P\mathbf{t}$ can be found by optimizing \tilde{f} along all the directions \mathbf{p}_j

$$\min_{\mathbf{t}} \quad \tilde{f}(\mathbf{w} + \sum_{j=1}^m t_j \mathbf{p}_j \mid \mathbf{w}) = (P^\top \nabla f(\mathbf{w}))^\top \mathbf{t} + \frac{1}{2} \mathbf{t}^\top (P^\top \nabla^2 f(\mathbf{w}) P) \mathbf{t}. \quad (6)$$

Problem (6) can be solved by the following linear system.

$$(P^\top \nabla^2 f(\mathbf{w}) P) \mathbf{t} = -P^\top \nabla f(\mathbf{w}).$$

Say we have stored m previous directions in $P \in \mathbb{R}^{m \times n}$, where $m \ll n$. Problem (6) has m variables but if instead we consider the Newton direction by solving the problem in (5), the number of variables is n . Therefore, solving (6) should be easier than finding a Newton direction if the construction of $P \nabla^2 f(\mathbf{w}) P^\top$ is not too expensive. In Section 3, we will show that when the Hessian is structured, the subproblem (6) can be solved easily.

From the above concept, we propose the common-directions method. In each iteration of this method, given the orthonormal basis P of the span of the past directions, we solve the subproblem (6) with respect to \mathbf{t} to obtain the update direction $P\mathbf{t}$. A line search along $P\mathbf{t}$ is performed to guarantee the strict decrease of the function value. After updating the iterate \mathbf{w} , a new direction $-\nabla f(\mathbf{w})$ is then taken to enlarge the orthonormal basis P by Gram-Schmidt orthogonalization. A description of the common-directions method is given in Algorithm 1, and the line-search sub-routine is described in Algorithm 2. Note that other choices of the directions being added to P are also possible, but for the ease of description and analysis, we confine our choice to the current steepest descent direction.

We define the following quantities before stating and proving the theorems. We denote the condition number of f by

$$\kappa \equiv \frac{\rho}{\sigma},$$

and define an ϵ -accurate solution of f to be a point \mathbf{w} such that

$$f(\mathbf{w}) - f(\mathbf{w}^*) \leq \epsilon,$$

where \mathbf{w}^* is the optimal solution of f . Note that under Assumption 2.1, problem (1) possesses a unique optimal solution. The following theorems present theoretical convergence of the proposed method. In particular, we show that the line search sub-routine has finite termination, and our method is globally linearly convergent and locally quadratically convergent.

Theorem 2.3 *Under Assumption 2.1, every time being called by Algorithm 1, Algorithm 2 terminates within $\lceil \log_\beta(\beta\sigma/(\rho + \lambda)) \rceil$ steps, and Algorithm 1 converges Q-linearly with iteration complexity $O(\kappa^3 \log(1/\epsilon))$ to achieve an ϵ -accurate solution.*

Theorem 2.4 *Under Assumptions 2.1 and 2.2, Algorithm 1 converges quadratically.*

The proofs are presented in Appendix A.

3. Application on Regularized Empirical Risk Minimization Problems

For most optimization problems, computing the $n \times n$ Hessian matrix, or its inverse, whose cost is $O(n^3)$, is much more expensive than obtaining the function value and the n -dimensional gradient vector. To tackle this problem, different approaches have been considered, and in some cases the problem structure can be utilized to reduce the cost. We are interested in solving machine learning problems, and thus we will focus our discussion on differentiable regularized ERM problems. More specifically, given training feature-label pairs (\mathbf{x}_i, y_i) , $i = 1, \dots, l$, with $\mathbf{x}_i \in \mathbb{R}^n$, $y_i \in \mathbb{R}$, and a parameter $C > 0$, we consider L2-regularized ERM problems

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}), \quad \text{where } f(\mathbf{w}) \equiv \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^l \xi(y_i; \mathbf{w}^\top \mathbf{x}_i), \quad (7)$$

and ξ is a loss function convex in $\mathbf{w}^\top \mathbf{x}_i$. The first term in (7) makes f be σ strongly convex with $\sigma \geq 1$, and the choice of loss functions ensures that $\nabla f(\mathbf{w})$ is Lipschitz continuous, though the corresponding constant ρ may be unknown. Therefore, Assumption 2.1 is satisfied.

For problem (7), by defining

$$X \equiv \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_l^\top \end{pmatrix},$$

we can write the gradient as

$$\nabla f(\mathbf{w}) = \mathbf{w} + X^\top \mathbf{v}_{\mathbf{w}}, \quad \text{where } (\mathbf{v}_{\mathbf{w}})_i = C \frac{\partial}{\partial z} \xi(y_i; z) |_{z=\mathbf{w}^\top \mathbf{x}_i}, i = 1, \dots, l, \quad (8)$$

and if $f(\mathbf{w})$ is twice-differentiable, the Hessian is

$$\nabla^2 f(\mathbf{w}) \equiv I + X^\top D_{\mathbf{w}} X, \quad (9)$$

where I is the identity matrix and $D_{\mathbf{w}}$ is a diagonal matrix with

$$(D_{\mathbf{w}})_{ii} \equiv C \frac{\partial^2}{\partial z^2} \xi(y_i; z) |_{z=\mathbf{w}^\top \mathbf{x}_i}, i = 1, \dots, l. \quad (10)$$

The structure of (9) allows efficient computations to use the Hessian in various ways. Our method utilizes it to incorporate the Hessian-vector products with low cost as discussed in the following sections.

Some problems of the form (7), such as logistic regression and ridge regression, also satisfy Assumption 2.2. The loss of ridge regression is

$$\xi(y_i; \mathbf{w}^\top \mathbf{x}_i) = (y_i - \mathbf{w}^\top \mathbf{x}_i)^2.$$

It is clearly twice-differentiable with the diagonal matrix $D_{\mathbf{w}}$ defined in (10) being the identity matrix. Clearly, this loss satisfies Assumption 2.2 with any $M \geq 0$ as the Hessian is a constant. The loss of logistic regression is

$$\xi(y_i; \mathbf{w}^\top \mathbf{x}_i) = \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)),$$

where $y_i \in \{-1, 1\}$. Because its third derivative is upper-bounded, by applying first-order Taylor expansion on the Hessian, one can see that Assumption 2.2 holds. On the other hand, problems that satisfy Assumption 2.1 only, like squared-hinge loss support vector classification/regression (Boser et al., 1992; Vapnik, 1995), can also be solved by our algorithm, although in this case we can only expect linear convergence.²

3.1 Subproblem of Common Directions

At the k -th iteration, given the matrix $P_k = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m] \in \mathbb{R}^{n \times m}$ with orthogonal columns, we consider solving (6) along $P_k \mathbf{t}$

$$\min_{\mathbf{t} \in \mathbb{R}^m} \tilde{f}(\mathbf{w}_k + P_k \mathbf{t} \mid \mathbf{w}_k). \quad (11)$$

The gradient and the Hessian of (11) with respect to the variable \mathbf{t} at $\mathbf{t} = \mathbf{0}$ are respectively

$$\nabla_{\mathbf{t}} \tilde{f}(\mathbf{w}_k + P_k \mathbf{t}) \mid_{\mathbf{t}=\mathbf{0}} = P_k^\top \nabla f(\mathbf{w}_k) \quad (12)$$

and

$$\nabla_{\mathbf{t}, \mathbf{t}}^2 \tilde{f}(\mathbf{w}_k + P_k \mathbf{t}) \mid_{\mathbf{t}=\mathbf{0}} = P_k^\top \nabla^2 f(\mathbf{w}_k) P_k = P_k^\top (I + X^\top D_{\mathbf{w}_k} X) P_k = I + P_k^\top X^\top D_{\mathbf{w}_k} X P_k. \quad (13)$$

Observe the term $X P_k$. The trick to make the computation efficient is to store this term in memory to avoid repetitive computations of $X \mathbf{p}_j$ for all previous directions \mathbf{p}_j . If $X P_k$, $\nabla f(\mathbf{w}_k)$, and $D_{\mathbf{w}_k}$ are available, the construction of (12)-(13) involves the following cost.

- $P_k^\top \nabla f(\mathbf{w}_k)$: $O(mn)$,
- $(X P_k)^\top D_{\mathbf{w}_k} (X P_k)$: $O(lm^2)$.

The optimal solution \mathbf{t} of (11) can be obtained by solving the following linear system

$$(I + P_k^\top X^\top D_{\mathbf{w}_k} X P_k) \mathbf{t} = -P_k^\top \nabla f(\mathbf{w}_k) \quad (14)$$

in $O(m^3)$ cost. Note that the Hessian with respect to \mathbf{t} is always invertible because the term I in (13) ensures its positive definiteness.

Let \mathbf{t}_k be the solution of (14). We show that to obtain necessary information for the next iteration, it is sufficient to maintain $X P_{k+1}$ and $X \mathbf{w}_{k+1}$. To have the new function value, from (7) we mainly need to compute $X \mathbf{w}_{k+1}$. If $X \mathbf{w}_k$ is available and the full direction $(X P_k) \mathbf{t}_k$ is taken, by

$$X \mathbf{w}_{k+1} = X \mathbf{w}_k + \theta(X P_k) \mathbf{t}_k, \quad (15)$$

with $\theta = 1$, in $O(lm)$ cost $X \mathbf{w}_{k+1}$ is obtained. In practice, a step size θ is decided by line search, which involves some minor cost as shown in Section 3.2. The vector $X \mathbf{w}_{k+1}$ is also used to construct $\mathbf{v}_{\mathbf{w}_{k+1}}$ needed for gradient calculation. Then from (8),

$$\nabla f(\mathbf{w}_{k+1}) = \mathbf{w}_{k+1} + X^\top \mathbf{v}_{\mathbf{w}_{k+1}}$$

2. If the objective function is not twice differentiable, we can used the generalized Hessian (Mangasarian, 2002) to replace the Hessian in the Newton step, as Assumption 2.1 guarantees the existence of generalized Hessian.

takes $O(\#\text{nnz})$ cost, where $\#\text{nnz}$ is the number of non-zero elements in X . We then apply the Gram-Schmidt procedure to check if a new column should be added to P_{k+1} . In practice this procedure is computed by

$$\mathbf{p}_{m+1} = (I - P_k P_k^T) \nabla f(\mathbf{w}_{k+1}) = \nabla f(\mathbf{w}_{k+1}) - P_k (P_k^\top \nabla f(\mathbf{w}_{k+1})).$$

The cost is $O(mn)$. If $\mathbf{p}_{m+1} \neq \mathbf{0}$, to maintain XP_{m+1} after \mathbf{p}_{m+1} is added, we must compute XP_{m+1} , which costs $O(\#\text{nnz})$. Details of the algorithm is shown in Algorithm 3.

In summary, the cost per iteration is

$$O(lm^2 + mn + m^3 + \#\text{nnz}). \quad (16)$$

Because usually we have $n \gg m$ and $l \gg m$ from the first iteration till reaching a good enough solution, in general $\#\text{nnz}$ is still the dominant term in (16). Therefore, each of our iterations does not cost much more than a typical iteration in an optimization method that often involves the $O(\#\text{nnz})$ cost.

Further, for the same set of directions P_k , we could solve the subproblem (11) and update \mathbf{w}_k multiple times in order to approximately solve (2). Note that the update of \mathbf{w}_k does not need to go through the data matrix X . The reason is that the D_w matrix and the \mathbf{v}_w vector rely on \mathbf{w}_k only through the inner products $\mathbf{w}_k^\top \mathbf{x}_i$, and this is a linear combination of XP_k , and thus we can compute it without additional inner product computations with the data matrix X . This idea of solving (11) multiple times will be exploited in Section 4 to achieve the optimal linear convergence rate for first-order methods.

3.2 Line Search on Common Directions

Once the solution \mathbf{t} of (14) has been obtained, in $O(mn)$ cost we get the descent direction

$$\mathbf{d} = P_k \mathbf{t}. \quad (17)$$

Next, we conduct line search in Algorithm 2 for the step size θ in

$$f(\mathbf{w}_k + \theta \mathbf{d}) = \frac{1}{2} \|\mathbf{w} + \theta \mathbf{d}\|^2 + C \sum_{i=1}^l \xi(y_i; \mathbf{w}_k^\top \mathbf{x}_i + \theta \mathbf{d}^\top \mathbf{x}_i).$$

The sufficient decrease condition is

$$C \sum_{i=1}^l \xi(y_i; \mathbf{w}_k^\top \mathbf{x}_i) - \theta \mathbf{w}^\top \mathbf{d} - \frac{\theta^2}{2} \|\mathbf{d}\|^2 - C \sum_{i=1}^l \xi(y_i; \mathbf{w}_k^\top \mathbf{x}_i + \theta \mathbf{d}^\top \mathbf{x}_i) < \frac{\lambda}{2} \theta^2 \|\mathbf{d}\|^2.$$

Because XP_k is maintained, after obtaining $X\mathbf{d}$ in $O(lm)$ time by

$$X\mathbf{d} = (XP_k)\mathbf{t}, \quad (18)$$

and calculating $\mathbf{w}^\top \mathbf{d}$, $\mathbf{w}^\top \mathbf{w}$, $\mathbf{d}^\top \mathbf{d}$ in $O(n)$ time, each line search step would only cost $O(l)$ because $\mathbf{w}_k^\top \mathbf{x}_i$ is already available in our calculation for the gradient and the Hessian. Thus, we could perform line search in

$$O(lm + mn + l \times \#(\text{line search steps})) \text{ time}. \quad (19)$$

By Theorem 2.3, $\#(\text{line search})$ is bounded by a constant, and thus the complexity does not grow infinitely. After finishing the line search procedure, we already have $\mathbf{w}_k^\top \mathbf{x}_i + \theta \mathbf{d}^\top \mathbf{x}_i$ as the next $\mathbf{w}_{k+1}^\top \mathbf{x}_i$ as shown in (15).

Algorithm 3: The common-directions method for solving the ERM problem (7)

```

Given  $\mathbf{w}_0$ , compute  $\mathbf{z} = X\mathbf{w}_0$ ;
Use  $\mathbf{z}$  to calculate  $\mathbf{v}_{\mathbf{w}_0}$  and  $D_{\mathbf{w}_0}$ ;
Compute  $\nabla f(\mathbf{w}_0) = \mathbf{w}_0 + X^\top \mathbf{v}_{\mathbf{w}_0}$  from (8);
 $P = [\nabla f(\mathbf{w}_0)/\|\nabla f(\mathbf{w}_0)\|], U = XP$ ;
for  $k=0,1,\dots$  do
    if  $\nabla f(\mathbf{w}_k) = \mathbf{0}$ , or some stopping condition is satisfied then
        | stop;
    end
    Obtain  $\mathbf{t}$  by solving (14) with  $\mathbf{t} = [I + U^\top D_{\mathbf{w}_k} U]^{-1}[-P^\top \nabla f(\mathbf{w}_k)]$ ;
    Compute  $\hat{\mathbf{z}} = U\mathbf{t}$ ;
    Backtracking line search on  $f(\mathbf{w}_k + \theta P\mathbf{t})$  by Algorithm 2 to obtain  $\theta$  using the
        values of  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  to compute  $\xi(y_i; \mathbf{z} + \theta\hat{\mathbf{z}})$ ;
     $\mathbf{w}_{k+1} = \mathbf{w}_k + \theta P\mathbf{t}$ ;
     $\mathbf{z} = \mathbf{z} + \theta\hat{\mathbf{z}}$ ;
    Use  $\mathbf{z}$  to calculate  $\mathbf{v}_{\mathbf{w}_{k+1}}$  and  $D_{\mathbf{w}_{k+1}}$ ;
    Compute  $\nabla f(\mathbf{w}_{k+1}) = \mathbf{w}_{k+1} + X^\top \mathbf{v}_{\mathbf{w}_{k+1}}$  from (8);
    Let  $\mathbf{p} = \nabla f(\mathbf{w}_{k+1}) - P(P^\top \nabla f(\mathbf{w}_{k+1}))$ ;
    if  $\mathbf{p} \neq \mathbf{0}$  then
        |  $P = [P, \mathbf{p}/\|\mathbf{p}\|], U = [U, X\mathbf{p}/\|\mathbf{p}\|]$ ;
    end
end

```

3.3 Total Cost of the Algorithm

Based on the analysis in Sections 3.1 and 3.2, the total cost is

$$O(lm^2 + mn + m^3 + \#\text{nnz} + l \times \#(\text{line search steps})) \times \text{iterations}. \quad (20)$$

In Section 5, we will compare the cost with that of some existing approaches.

4. Common-directions Method with Multiple Inner Iterations

In Section 2, we considered solving the second-order Taylor approximation of $f(\mathbf{w}_k + P_k \mathbf{t})$ with respect to the coefficients \mathbf{t} to get a good combination along the directions of P_k . Now we further extend the idea to directly minimize $f(\mathbf{w}_k + P_k \mathbf{t})$. We first describe the general framework, and then discuss how it can be applied to the special case of solving regularized ERM problems.

4.1 General Framework

We can consider the following convex optimization problem

$$\min_{\mathbf{t}} f(\mathbf{w}_k + P_k \mathbf{t}), \quad (21)$$

and apply any optimization algorithm to solve it. When applying iterative algorithms, we terminate the optimization procedure after the solution is considered good enough.

Therefore, usually (21) is only solved approximately. In particular, the method discussed in Section 2 can be viewed as using only one Newton iteration to loosely solve (21). Because of two levels of iterations, we refer to the process of approximately solving (21) as an outer iteration, while each step in the process as an inner iteration.

Among the different methods to solve (21), we apply multiple inner Newton steps on the variable \mathbf{t} . That is, we repeatedly minimize the second-order Taylor expansion of (21).

$$\min_{\Delta \mathbf{t} \in \mathbb{R}^m} \tilde{f}(\mathbf{w}_k + P_k \mathbf{t} + P_k \Delta \mathbf{t} \mid \mathbf{w}_k + P_k \mathbf{t}).$$

Because

$$\begin{aligned} & \tilde{f}(\mathbf{w}_k + P_k \mathbf{t} + P_k \Delta \mathbf{t} \mid \mathbf{w}_k + P_k \mathbf{t}) \\ &= f(\mathbf{w}_k + P_k \mathbf{t}) + (P_k^\top \nabla f(\mathbf{w}_k + P_k \mathbf{t}))^\top \Delta \mathbf{t} + \frac{1}{2} \Delta \mathbf{t}^\top P_k^\top \nabla f(\mathbf{w}_k + P_k \mathbf{t}) P_k \Delta \mathbf{t}, \end{aligned}$$

this amounts to solving the following linear system for $\Delta \mathbf{t}$.

$$(P_k^\top \nabla^2 f(\mathbf{w}_k + P_k \mathbf{t}) P_k) \Delta \mathbf{t} = -P_k^\top \nabla f(\mathbf{w}_k + P_k \mathbf{t}). \quad (22)$$

Each time the solution of the linear system will be our update direction, and then we will conduct backtracking line search to ensure the sufficient function decrease. The detailed algorithm is described in Algorithm 4, where we can see that $\mathbf{w}_k + P_k \mathbf{t}$ is the iterate updated in the inner iterations.

We now show that this general framework has better global convergence rate than the one discussed in Section 2.

Theorem 4.1 *Under Assumption 2.1 and a proper inner stopping condition, given any $\epsilon > 0$, Algorithm 4 converges R-linearly and obtains an ϵ -accurate solution in $O(\sqrt{\kappa} \log(1/\epsilon))$ iterations.*

The proof is in Appendix B.

It has been shown in Nesterov (2003) that $O(\sqrt{\kappa} \log(1/\epsilon))$ is the optimal iteration complexity for methods whose directions are obtained from the span of $\nabla f(\mathbf{w}_0), \nabla f(\mathbf{w}_1), \dots$. The inner stopping condition used in proving Theorem 4.1 is not easy to implement in practice. It is possible to devise a practical one to ensure that this optimal linear convergence rate can be guaranteed. However, as we will see in the empirical results in Section 6, the setting of using a single inner iteration in approximately solving (21) tends to outperform that of using multiple inner iterations in terms of overall training time. In light of this understanding, we do not further explore practical inner stopping conditions.

Another interesting property of Algorithm 4 is that if the optimal solution of (21) is obtained at each outer iteration, then the algorithm can reach the optimum of (7) in finite steps even if neither Assumption 2.1 nor 2.2 is satisfied.

Theorem 4.2 *If at each iteration of Algorithm 4, the optimal solution of (21) is obtained, the optimum of (1) is reached within n iterations.*

The proof is in Appendix C.

These two theorems show that our algorithm can be treated as an analogy of the linear conjugate gradient method to general nonlinear optimization problems. See more details in Section 5.1

Algorithm 4: Common-directions method with multiple inner iterations

```

Given  $\mathbf{w}_0$ , compute  $\nabla f(\mathbf{w}_0)$ ;
 $P = [\nabla f(\mathbf{w}_0)/\|\nabla f(\mathbf{w}_0)\|]$ ;
for  $k=0,1,\dots$  do
    if  $\nabla f(\mathbf{w}_k) = \mathbf{0}$ , or some stopping condition is satisfied then
        | stop;
    end
     $\mathbf{w} = \mathbf{w}_k$ ;
    while an inner stopping condition is not satisfied do // inner loop
        | Obtain  $\Delta\mathbf{t}$  by solving (22) with
            
$$(P^\top \nabla^2 f(\mathbf{w}) P) \Delta\mathbf{t} = -P^\top \nabla f(\mathbf{w});$$

        | Backtracking line search on  $f(\mathbf{w} + \theta P \Delta\mathbf{t})$  by Algorithm 2 to obtain  $\theta$ ;
        |  $\mathbf{w} = \mathbf{w} + \theta P \Delta\mathbf{t}$ ;
    end
     $\mathbf{w}_{k+1} = \mathbf{w}$ ;
    Compute  $\nabla f(\mathbf{w}_{k+1})$ ;
    Let  $\mathbf{p} = \nabla f(\mathbf{w}_{k+1}) - P(P^\top \nabla f(\mathbf{w}_{k+1}))$ ;
    if  $\mathbf{p} \neq \mathbf{0}$  then
        |  $P = [P, \mathbf{p}/\|\mathbf{p}\|]$ ;
    end
end

```

4.2 Application on Regularized ERM Problems

For general problems, Algorithm 4 may not be beneficial because the cost of one inner iteration of solving (21) is similar to that of conducting an outer iteration of Algorithm 1, while (1) rather than (21) is what we actually would like to solve. However, for solving the regularized ERM problems (7), we show that one inner iteration of the while-loop in Algorithm 5 can be significantly cheaper than one outer iteration of Algorithm 3. We note that in the complexity analysis (16) for Algorithm 3, the dominant term $O(\#\text{nnz})$ is unavoidable for calculating $\nabla f(\mathbf{w}_{k+1})$. A crucial difference here is that this $O(\#\text{nnz})$ term is not needed. In Algorithm 3, the whole $\nabla f(\mathbf{w}_{k+1})$ vector is used for the Gram-Schmidt procedure in checking if P_k should be augmented. Instead, here P_k remains the same throughout inner iterations and all we need is $P_k^\top \nabla f(\mathbf{w}_k + P_k \mathbf{t})$ in (22). When we apply Algorithm 4 to regularized ERM problems, the linear system (22) has the following form.

$$(I + P_k^\top X^\top D_{\mathbf{w}_k + P_k \mathbf{t}} X P_k) \Delta\mathbf{t} = -(P_k^\top (\mathbf{w}_k + P_k \mathbf{t}) + P_k^\top X^\top \mathbf{v}_{\mathbf{w}_k + P_k \mathbf{t}}). \quad (23)$$

If

$$X P_k \text{ and } X(\mathbf{w}_k + P_k \mathbf{t}), \quad (24)$$

are maintained, the right-hand side can be easily calculated by

- $\mathbf{v}_{\mathbf{w}_k + P_k \mathbf{t}} : O(l)$,

- $P_k^\top(\mathbf{w}_k + P_k \mathbf{t}) : O(mn),$
- $(P_k^\top X^\top) \mathbf{v}_{\mathbf{w}_k + P_k \mathbf{t}} : O(lm).$

Interestingly, for other operations in Algorithm 3, it has been shown that the two values in (24) can be easily maintained. Here we follow the same setting. A detailed procedure is in Algorithm 5. In summary, the cost per inner iteration is

$$O(lm^2 + mn + m^3). \quad (25)$$

In compared with the cost per outer iteration in Algorithm 3, here we do not have the $\#\text{nnz}$ term. The main difference comes from that in (23), we use the available XP_k to calculate

$$(P_k^\top X^\top) \mathbf{v}_{\mathbf{w}_k + P_k \mathbf{t}},$$

but in an outer iteration of Algorithm 3, $\nabla f(\mathbf{w}_{k+1})$ is calculated by

$$X^\top \mathbf{v}_{\mathbf{w}_{k+1}},$$

which costs $O(\#\text{nnz})$. Therefore, because in general $\#\text{nnz}$ is larger than (25), an inner iteration here is cheaper than an outer one in Algorithm 3.

5. Related Works

We compare our algorithm with linear conjugate gradient method for solving a linear system, as well as popular nonlinear optimization methods. The methods being discussed include those considering a linear combination of multiple directions to form the update direction, and those utilizing the Hessian matrix in various ways.

5.1 Linear Conjugate Gradient Method

Linear conjugate gradient (CG) method iteratively solves the following linear system.

$$A\mathbf{x} = \mathbf{b}, \quad (26)$$

where $A \in \mathbb{R}^{n \times n}$ is symmetric and positive-definite, and $\mathbf{b} \in \mathbb{R}^n$. Alternatively, CG can be viewed as an optimization method for the following problem.

$$\min_{\mathbf{x}} \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|^2. \quad (27)$$

The idea behind linear CG is that at each iteration, we generate a new update direction \mathbf{p}_k that is conjugate to the existing ones such that

$$\mathbf{p}_i^\top A \mathbf{p}_j = 0, \quad \forall i \neq j.$$

Then an exact line search is conducted to decide the next iterate. At each iteration, a matrix-vector product $A\mathbf{v}$ for some vector \mathbf{v} is conducted to generate the desired step, and this matrix-vector product is the computational dominant part at each CG iteration.

Algorithm 5: Common-directions method with multiple inner iterations for the ERM problem (7)

```

Given  $\mathbf{w}_0$ , compute  $\mathbf{z} = X\mathbf{w}_0$ ;
Use  $\mathbf{z}$  to compute  $\mathbf{v}_{\mathbf{w}_0}$  and  $D_{\mathbf{w}_0}$  in (8) and (10);
Compute  $\nabla f(\mathbf{w}_0) = \mathbf{w}_0 + X^\top \mathbf{v}_{\mathbf{w}_0}$  from (8);
 $P = [\nabla f(\mathbf{w}_0)/\|\nabla f(\mathbf{w}_0)\|]$ ,  $U = XP$ ;
for  $k=0,1,\dots$  do
    if  $\nabla f(\mathbf{w}_k) = \mathbf{0}$ , or some stopping condition is satisfied then
        | stop;
    end
     $\mathbf{w} = \mathbf{w}_k, \mathbf{v}_{\mathbf{w}} = \mathbf{v}_{\mathbf{w}_k}, D_{\mathbf{w}} = D_{\mathbf{w}_k}$ ;
    while an inner stopping condition is not satisfied do // inner loop
        | Obtain  $\Delta t$  by solving (23) with
            
$$\Delta t = (I + U^\top D_{\mathbf{w}} U)^{-1}(-P^\top \mathbf{w} + U^\top \mathbf{v}_{\mathbf{w}});$$

        | Compute  $\hat{\mathbf{z}} = U\Delta t$ ;
        | Backtracking line search on  $f(\mathbf{w} + \theta P \Delta t)$  by Algorithm 2 to obtain  $\theta$ , using
            the values of  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  to compute  $\xi(y_i; \mathbf{z} + \theta \hat{\mathbf{z}})$ ;
        |  $\mathbf{w} = \mathbf{w} + \theta P \Delta t$ ;
        |  $\mathbf{z} = \mathbf{z} + \theta (XP) \Delta t$ ;
        | Use  $\mathbf{z}$  to compute  $\mathbf{v}_{\mathbf{w}}$  and  $D_{\mathbf{w}}$ ;
    end
     $\mathbf{w}_{k+1} = \mathbf{w}, \mathbf{v}_{\mathbf{w}_{k+1}} = \mathbf{v}_{\mathbf{w}}, D_{\mathbf{w}_{k+1}} = D_{\mathbf{w}}$ ;
    Compute  $\nabla f(\mathbf{w}_{k+1}) = \mathbf{w}_{k+1} + X^\top \mathbf{v}_{\mathbf{w}_{k+1}}$  from (8);
     $\mathbf{p} = \nabla f(\mathbf{w}_{k+1}) - P(P^\top \nabla f(\mathbf{w}_{k+1}))$ ;
    if  $\mathbf{p} \neq \mathbf{0}$  then
        |  $P = [P_k, \mathbf{p}/\|\mathbf{p}\|]$ ,  $U = [U, X\mathbf{p}/\|\mathbf{p}\|]$ ;
    end
end

```

Therefore, the cost per iteration of CG is simply the cost of a matrix-vector product and several $O(n)$ operations on the vectors.

There are two nice properties of linear CG. First, when viewed as a solver for (26), it obtains the exact solution in n steps. Second, when viewed as a solver for (27), its iteration complexity for an ϵ -accurate solution is $O(\sqrt{\kappa} \log(1/\epsilon))$. See, for example, Chapter 5.1 of Nocedal and Wright (2006). Because the method proposed in Section 4 possesses the same two properties in Theorems 4.1 and 4.2, it can be viewed as an extension of linear CG to general nonlinear convex optimization.

5.2 First-order Methods and Momentum

For nonlinear convex optimization, the simplest first-order method is gradient descent. In gradient descent, the update direction is always the negative gradient $-\nabla f(\mathbf{w})$. It has been

shown in many standard convex optimization books such as [Boyd and Vandenberghe \(2004\)](#); [Nocedal and Wright \(2006\)](#); [Nesterov \(2003\)](#) that with a proper choice of the step size, to obtain an ϵ -accurate solution for a given $\epsilon > 0$, gradient descent requires $O(\kappa \log(1/\epsilon))$ iterations. However, this is far from the aforementioned optimal convergence rate $O(\sqrt{\kappa} \log(1/\epsilon))$ for methods that have their update directions in the span of all previous gradients.

Methods like the heavy-ball method ([Polyak, 1964](#)) and nonlinear CG try to improve the convergence by using the “momentum.” These methods possess update rules in the form

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) + \beta_k (\mathbf{w}_k - \mathbf{w}_{k-1}), \quad (28)$$

where α_k and β_k vary in different methods, and the term $\mathbf{w}_k - \mathbf{w}_{k-1}$ is the so-called momentum. By induction we can easily see that $\mathbf{w}_k - \mathbf{w}_{k-1}$ in (28) is a linear combination of $\nabla f(\mathbf{w}_0), \nabla f(\mathbf{w}_1), \dots, \nabla f(\mathbf{w}_{k-1})$, and therefore these methods fall in the category of picking the update direction from the span of the previous gradients. Nonlinear CG has many different variants that consider different strategies based on \mathbf{w}_k , \mathbf{w}_{k-1} , $\nabla f(\mathbf{w}_k)$, $\nabla f(\mathbf{w}_{k-1})$ and even $\nabla^2 f(\mathbf{w}_k)$, and all of them reduce to linear CG when the problem is (27). Interested readers are referred to [Hager and Zhang \(2006\)](#) and the references therein. When f is quadratic and strongly convex, it can be shown that both the heavy-ball method with proper parameters ([Lessard et al., 2016](#)) and nonlinear CG methods can achieve the optimal convergence rate for first-order methods,³ but for non-quadratic functions, the situation is different. Some variants of nonlinear CG can be shown to converge asymptotically, but the rate is unclear ([Hager and Zhang, 2006](#), Theorem 5.1). On the other hand, [Ghadimi et al. \(2014\)](#) showed that for problems satisfying Assumption 2.1 and twice-differentiability, the convergence rate of the heavy-ball method is the same as gradient descent.

Nesterov’s accelerated gradient ([Nesterov, 1983](#)) considers two sequences $\{\mathbf{w}_k\}$ and $\{\mathbf{s}_k\}$ with the following update rules.

$$\begin{aligned} \mathbf{w}_{k+1} &\equiv \mathbf{s}_k - \frac{1}{\rho} \nabla f(\mathbf{s}_k), \\ \mathbf{s}_{k+1} &= (1 + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}) \mathbf{w}_{k+1} - \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \mathbf{w}_k. \end{aligned}$$

Interestingly, this method uses only the previous gradients in a deterministic way with κ to be the only data-dependent value, but it can guarantee the optimal iteration complexity of first-order methods for problems satisfying Assumption 2.1. A problem in this method is that we need to know the values ρ and σ in advance, and these values may not be available. To solve this problem, a variant of Nesterov’s accelerated gradient is proposed in [Nesterov \(2013\)](#). This algorithm estimates these parameters at each iteration via an iterative procedure similar to backtracking line search, and achieves the same iteration complexity with little additional cost. Note that for all these first-order methods, the main bottleneck is the computation of the gradients, and the cost is proportional to the number of different gradients computed per iteration.

Our method can be viewed as a generalization of these approaches, in the sense that the same search space for the update direction is considered. The major difference is

3. Section 1 of [Hager and Zhang \(2006\)](#) suggests that nonlinear CG methods reduce to linear CG when f is quadratic, and the convergence rate follows what we described in Section 5.1.

that these momentum methods are restricted to some deterministic or pre-specified ways of combination, while our method considers combining multiple directions with adaptive weights according to the data that can and will change over iterations. Another difference is that we utilize higher-order derivatives to decide the combination coefficient, so quadratic convergence in Theorem 2.4 can be obtained.

5.3 Methods Utilizing Multiple Directions

The cutting plane method (Kelley, 1960) and the bundle method (Teo et al., 2010) use all gradients obtained from the first iteration on. They differ from our method from using the previous gradients to construct a piecewise linear function to approximate (the hard-to-compute part of) the actual function.⁴ Therefore the purpose of using previous gradients in these methods and that of ours are different. An advantage of the cutting plane method and the bundle method is the ability to deal with nonsmooth problems by using subgradients. When being applied to problems satisfying Assumption 2.1, the bundle method converges linearly (Teo et al., 2010). However, optimizing a piecewise linear function with a growing number of pieces is more expensive than dealing with one smooth function. Our method considers easier-to-solve subproblems at each iteration. Moreover, our method converges quadratically while the bundle method and the cutting plane method do not have this property.

Quasi-Newton methods also consider previous momentums and gradients, but the purpose is to approximate the real Hessian matrix via these vectors in order to achieve superlinear convergence. Among different quasi-Newton methods, the BFGS method (Dennis and Moré, 1977) and its limited-memory variant, L-BFGS (Liu and Nocedal, 1989), are the most widely used ones. By applying a proper line search procedure, the two methods have global linear convergence for problems satisfying Assumption 2.1. Note that because older information are discarded, L-BFGS does not have superlinear convergence, but it has the advantage of less spatial cost and is more suitable for large-scale problems.

At the k -th iteration for any $k \geq 1$, given the current iterate \mathbf{w}_k , BFGS solves the following quadratic problem.

$$\min_{\mathbf{p}_k} \nabla f(\mathbf{w}_k)^\top \mathbf{p}_k + \frac{1}{2} (\mathbf{p}_k)^\top B_k \mathbf{p}_k, \quad (29)$$

where B_k is a positive definite symmetric approximation of $\nabla^2 f(\mathbf{w}_k)$ such that

$$B_k^{-1} = (I - \mu_{k-1} \mathbf{u}_{k-1} \mathbf{s}_{k-1}^\top)^\top B_{k-1}^{-1} (I - \mu_{k-1} \mathbf{u}_{k-1} \mathbf{s}_{k-1}^\top) + \mu_{k-1} \mathbf{s}_{k-1} \mathbf{s}_{k-1}^\top, \quad (30)$$

with

$$\mu_{k-1} \equiv \frac{1}{\mathbf{u}_{k-1}^\top \mathbf{s}_{k-1}}, \quad \mathbf{s}_{k-1} = \mathbf{w}_k - \mathbf{w}_{k-1}, \quad \mathbf{u}_{k-1} = \nabla f(\mathbf{w}_k) - \nabla f(\mathbf{w}_{k-1}).$$

Note that although B_0 can be decided arbitrarily, the most common choice is

$$B_0 = \beta I \quad (31)$$

4. It is known in functional analysis that any convex function can be approximated by piecewise linear functions to any precision.

for some $\beta > 0$. The solution of (29) is clearly

$$\mathbf{p}_k = -B_k^{-1}\nabla f(\mathbf{w}_k). \quad (32)$$

If (31) is used, this operation can be conducted cheaply with only $O(kn)$ cost by a sequence of vector-vector operations obtained from expanding (30).⁵ After deciding \mathbf{p}_k , a line search procedure is conducted to find a suitable step size θ_k to ensure that the sequence $\{\mathbf{w}_k\}$ converges to the optimal solution of f .

With the choice of (31), by induction the vector

$$\mathbf{s}_k = \alpha_{k+1}\mathbf{p}_k \quad (33)$$

is a linear combination of $\nabla f(\mathbf{w}_0), \dots, \nabla f(\mathbf{w}_k)$ for all k .⁶ Thus we can see that $\mathbf{p}_k \in \text{span}\{\nabla f(\mathbf{w}_0), \dots, \nabla f(\mathbf{w}_k)\}$ and thus the BFGS algorithm only utilizes the previous gradients, but is able to achieve superlinear convergence. Interestingly, our algorithm utilizes the same information as BFGS, but because we also consider the Hessian to search for the best update direction in $\text{span}\{\nabla f(\mathbf{w}_0), \dots, \nabla f(\mathbf{w}_k)\}$, our method is better than BFGS for being able to achieve quadratic convergence.

A major drawback of the BFGS algorithm is that it requires storing many vectors (or an $O(n^2)$ dense matrix) in memory. The memory consumption can be prohibitively expensive after iterations. This issue also happens in our method, but we notice that in practice, it is not severe. One reason is that the memory capacity is large enough to store the required information unless the number of iterations is high, and both our method and BFGS converge fast enough such that the memory capacity is not reached. For extremely large data, the memory issue may be faced, but this issue is less problematic in our method than in BFGS. The reason is that at each iteration of BFGS, two new vectors are added to memory, while in our method, at most one vector is added.

To handle the memory consumption problem of BFGS, its limited-memory variant, L-BFGS (Liu and Nocedal, 1989), is proposed. L-BFGS uses a different setting of B_k^{-1} . Given an integer $m > 0$, at the k -th iteration, L-BFGS only uses $\mathbf{s}_i, \mathbf{u}_i, i = k-m, k-m+1, \dots, k-1$ and $\nabla f(\mathbf{w}_k)$ to construct the approximate Hessian inverse. For L-BFGS, in the recursive computation of (30), the matrix B_{k-m}^{-1} is replaced by

$$\frac{1}{\mu_{k-m}}I.$$

As mentioned earlier, L-BFGS has only linear but not superlinear convergence. Nevertheless, it might outperform BFGS practically because the computation cost per iteration and the memory access amount are both lower. Moreover, it is shown in Bubeck et al. (2015) that L-BFGS outperforms those methods mentioned in Section 5.2 empirically, although the coefficient of its linear convergence is not clear.

-
- 5. If B_0 has neither sparsity nor some special structures, the cost of computing (32) at each iteration is $O(n^2)$ by directly maintaining (30), and it does not fall in the category of algorithms that only use gradients anymore because the vector $B_0^{-1}\nabla f(\mathbf{w})$ is considered.
 - 6. By continually expanding (32) using (31) and (30), we can see that \mathbf{p}_k is a linear combination of $\mathbf{s}_i, \mathbf{u}_i, i = 1, \dots, k-1$ and $\nabla f(\mathbf{w}_k)$. By induction, if all \mathbf{s}_i are linear combinations of previous gradients, then with the definition of \mathbf{u}_i , \mathbf{p}_k is also a linear combination of previous gradients. The induction is finished by that \mathbf{s}_k is a multiple of \mathbf{p}_k in (33).

5.4 Newton Method

At each iteration, given the current iterate \mathbf{w} , the Newton method obtains the update direction by solving

$$\min_{\mathbf{p} \in \mathbb{R}^n} \mathbf{p}^\top \nabla f(\mathbf{w}) + \frac{1}{2} \mathbf{p}^\top \nabla^2 f(\mathbf{w}) \mathbf{p}. \quad (34)$$

To obtain the solution of (34), the following n by n linear system is solved.

$$\nabla^2 f(\mathbf{w}) \mathbf{p} = -\nabla f(\mathbf{w}) \Rightarrow \mathbf{p} = -\nabla^2 f(\mathbf{w})^{-1} \nabla f(\mathbf{w}). \quad (35)$$

Then a line search procedure is conducted to ensure the sufficient function decrease. It is known that for twice-differentiable functions satisfying Assumption 2.2, the Newton method with line search has local quadratic convergence (Boyd and Vandenberghe, 2004; Nocedal and Wright, 2006). However, The computational cost per Newton iteration for directly solving the linear system (35) is as expensive as $O(n^3)$. Moreover, solving (35) precisely may not be beneficial when the current iterate is far from optimum. Therefore, truncated Newton methods that utilize, for example, the iterative linear CG method we discussed in Section 5.1 to approximately solve (35) is often used. In general, the computational cost per CG iteration is $O(n^2)$ for conducting the Hessian-vector product provided the Hessian is available. When we only want to solve the problem approximately, usually the required CG iteration is much smaller than n , making the computation of a truncated Newton step much cheaper than that of a full Newton step. It is known that truncated Newton steps combined with a line search procedure or trust region techniques retain local quadratic convergence for problems satisfying Assumption 2.2 and have superlinear convergence for problems satisfying Assumption 2.1, if the residual of the solution to (35) goes to zero as the number of truncated Newton iterations goes to infinity. See, for example, Lin and Moré (1999) and (Nocedal and Wright, 2006, Theorem 7.2).

As mentioned above, the cost of directly computing the Hessian-vector product is $O(n^2)$ both spatially and computationally if we store the Hessian matrix explicitly. However, as Keerthi and DeCoste (2005); Lin et al. (2008) proposed, for ERM problems in (7), if we calculate the product between Hessian and a given vector \mathbf{v} by

$$\nabla^2 f(\mathbf{w}) \mathbf{v} = (I + X^\top D_{\mathbf{w}} X) \mathbf{v} = \mathbf{v} + \left(X^\top (D_{\mathbf{w}}(X \mathbf{v})) \right), \quad (36)$$

this computation costs only $O(\#\text{nnz} + n)$ per CG iteration, where $\#\text{nnz}$ is at most ln but is usually much smaller for sparse data. Further, no storage space is needed for the operation in (36). On the other hand, explicitly forming the Hessian matrix in this case costs $O(\#\text{nnz}^2)$. Therefore, for regularized ERM problems, the CG approach without explicitly computing the Hessian matrix is more preferable. In this sense, truncated Newton methods costs

$$O((\#\text{nnz} + n) \times \text{CG iterations}) \quad (37)$$

per Newton iteration. A comparison between (37) and (20) shows that by a careful utilization of the problem structure of ERM, we have a cheaper cost of running one Newton iteration to approximately minimize a sub-problem over a subspace. The cost of our Newton iteration is only about that of one CG iteration here. We will examine the relative training speed between the truncated Newton method and our algorithm by experiments in Section 6.

Data	Training size (l)	Features (n)	Density (#nnz/(ln))
a9a	32,561	123	11.2757%
covtype	581,012	54	21.9979%
news20	19,996	1,355,191	0.0336%
url	2,396,130	3,231,961	0.0036%
epsilon	400,000	2000	100.0000%
webspam	350,000	16,609,143	0.0224%
yahoo-japan	140,963	832,026	0.0160%
yahoo-korea	368,444	3,053,939	0.0111%
rcv1t	677,399	47,236	0.1549%
KDD2010-b	19,264,097	29,890,095	0.0001%

Table 1: Statistics of data sets.

6. Experiments

We examine the empirical speed of our algorithm for training L2-regularized logistic regression models, in which $\xi(\cdot)$ in (7) is defined by

$$\xi(y; z) \equiv \log(1 + \exp(-yz))$$

with $y \in \{-1, 1\}$. We first compare the approach of conducting only one inner iteration discussed in Section 2 with that of conducting multiple inner iterations illustrated in Section 4. The relative efficiency between our methods and the methods reviewed in Section 5 is then examined.

6.1 Experiment Settings

We consider the data sets listed in Table 1 with three different choices of parameters $C = \{10^{-3}, 1, 10^3\}$ to examine the situation of different condition numbers. All data sets except yahoo-japan and yahoo-korea are publicly available.⁷ We use the relative distance to the optimal objective value, which is defined below, to compare different optimization approaches.

$$\left| \frac{f(\mathbf{w}) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \right|,$$

where \mathbf{w}^* is the optimal solution obtained by running our algorithm long enough. All experiments are conducted on a 64-bit machine with Intel Xeon 2.0 GHz CPU (E5-2620), 15MB cache, and 128GB memory.

We consider two different criteria for comparison. The first one is the empirical training time. This is the most important criterion in single-machine training for most users. We also consider the number of data passes, which is proportional to the number of iterations⁸ and is thus an indicator of iteration complexity. This criterion can also serve as a training time

7. <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>.

8. For the Newton method, this is the number of inner iterations in CG plus the number of Newton iterations.

estimate for scenarios that the cost is highly dependent to the number of passes through the data, such as distributed machine learning, or disk-level classification (Yu et al., 2012). In addition, this criterion excludes the effect of implementation differences.

6.2 Using Single or Multiple Inner Iterations

We compare the two variants of our algorithms, namely the approach of using only a single inner iteration and the variant of using multiple inner iterations, in Figures 1-6. Note that the result of KDD2010-b with $C = 10^3$ is not shown because the training time is too long. For the backtracking line search procedure in Algorithm 2, we set $\beta = 0.4, \lambda = 0.25$. The solution to the linear system (14) is obtained by a Cholesky factorization with partial pivoting. For the multiple inner iteration variant, we use the following heuristic stopping condition for the inner loop in Algorithm 4.

$$\|P_k \nabla f(\mathbf{w})\| \leq \min(0.1 \|\nabla f(\mathbf{w}_k)\|, \|\nabla f(\mathbf{w}_k)\|^2).$$

We can see that the two approaches perform quite similarly in both criteria, while the single inner iteration approach is slightly better in the training time. The figures of data passes indicate that using a single inner iteration has similar function value decrease per iteration to using multiple inner iterations, so the additional cost of multiple inner iterations sometimes results in slightly longer training time. Therefore, we use the single iteration variant in later experiments of comparing with existing algorithms.

6.3 Comparison Between Different Methods

We compare the following different algorithms for training logistic regression.

- CommDir: our common-directions method with single inner iteration.
- Line search truncated Newton method (NEWTON): at each Newton iteration, (34) is approximately solved by linear CG, until CG generates an approximate solution \mathbf{p} of (34) that satisfies

$$\|\nabla f(\mathbf{w}) - \nabla^2 f(\mathbf{w})\mathbf{p}\| \leq 0.1 \|\nabla f(\mathbf{w})\|.$$

We follow the discussion in Section 5 to efficiently compute the Hessian-vector products in $O(n + \#nnz)$ time. The implementation is modified from that in LIBLINEAR (Fan et al., 2008).

- L-BFGS: We implement the L-BFGS algorithm by considering different numbers of history states $m = 5, 10, 15, 20, 30$. In general $m = 30$ performs the best, so we report results of using this number of states.⁹
- BFGS: we use our L-BFGS implementation and set m to infinity to obtain a BFGS solver.

9. We also tried the implementation by the authors of Liu and Nocedal (1989) at <http://users.iems.northwestern.edu/~nocedal/lbfgs.html>, but found that it is not as efficient as ours. Note that their implementation does not utilize the efficient line search procedure discussed in Section 3.2.

- Nesterov’s accelerated gradient (AG): we consider the approach proposed by [Nesterov \(2013\)](#) that adaptively estimates the parameters σ and ρ by a procedure similar to backtracking line search. We take the initial estimation to be $\sigma = 1$ and $\rho = 1/C_l$.

For NEWTON, L-BFGS and BFGS that require a line search procedure to ensure the convergence, we follow our method to apply Algorithm 2 with the same parameters $\beta = 0.4$ and $\lambda = 0.25$ used in Section 6.2. We exploit the trick discussed in Section 3.2 to efficiently conduct this line search procedure. Note that for all methods except AG, evaluating $X^T \mathbf{v}_w$ in (8) for gradient at each iteration requires one data pass, and the line search procedure requires another one to maintain $X\mathbf{w}$. For NEWTON, each CG iteration requires one data pass. For AG, each inner iteration for adaptively estimating σ and ρ requires two data passes to compute related values. All algorithms are implemented in C++.

The results are shown in Figures 7-9 for running time and Figures 10-12 for number of data passes. Note that the result of KDD2010-b with $C = 10^3$ is not shown because the training time is too long. We can see that our method has the fewest data passes in all settings and all data sets. This observation suggests the usefulness of our method in the scenario that the data passes are expensive. The number of data passes is significantly reduced because we use the ideas in Section 3 to store the inner products XP_k to avoid redundant data accesses. In terms of training time, our method is among the fastest. The exceptions are **a9a** and **rcv1t** whose l are much larger than n , and the largest data set KDD2010-b, which from Table 1 is highly sparse. From the result of data passes, we see that our method requires fewer iterations than other methods, in the cost of the $O(lm^2 + mn)$ operations per iteration in addition to the $O(\#nnz)$ cost required by all methods to conduct one data pass; see (20). When $O(lm^2 + mn) \ll O(\#nnz)$, this additional cost is negligible and hence the advantage of fewer iterations dominates. However, when n is relatively small, or when the data set is highly sparse, the $O(lm^2 + mn)$ cost is significant in compared with $O(\#nnz)$, so our method may be slower for these data sets.

Note that the comparison between NEWTON and L-BFGS is inconsistent with that in [Lin et al. \(2008\)](#), because each line search iteration in the experiments of [Lin et al. \(2008\)](#) requires one data pass while we use the idea of caching inner products to accelerate the line search in our experiments. Moreover, in [Lin et al. \(2008\)](#), instead of line search, a trust region approach is considered for the truncated Newton method.

An interesting observation in our experimental results is the impressive performance of BFGS. Although BFGS has better convergence than L-BFGS, in practice people tend to use L-BFGS because BFGS consumes too much memory. Our experiments indicate that with larger memory and better computation power nowadays, BFGS becomes feasible and can outperform L-BFGS in some cases.

7. Discussion

As shown in our experiments, our method has a smaller number of data passes than all other methods. This property implies that our method has at least two potential applications for solving extremely large-scale ERM problems. The first one is disk-level classification when the data cannot fit into memory, and the second one is distributed machine learning. In both tasks, each data pass involves expensive operations like intensive disk I/O or between-

machine communication. With fewer data passes, our method is suitable for these tasks to tackle the bottleneck.

Each data pass in disk-level classification requires expensive disk I/O to access data because we cannot store all of them in memory. Existing methods for disk-level classification (Yu et al., 2012; Chang and Roth, 2011) consider solving the dual problem with techniques involving shuffling and partitioning data into blocks. Each block can be fully loaded into memory and these methods work on one block at a time. The random shuffling step gives better convergence, but requires lengthy pre-processing time because of heavy disk I/O. In contrast, our method does not require the random shuffling of data.

To reduce the number of data passes, Chang and Roth (2011) considered caching important instances in memory. Our method also uses ideas similar to caching, while the information being kept in memory is not some training instances, but the inner products XP_k . The number of these inner products is roughly the same as the number of iterations, which is usually not too large. Therefore, the memory consumption of our method tends to be smaller, and hence our method can cache more important information to reduce the data I/O.

For distributed machine learning that uses multiple machines, each data pass requires an expensive communication of a vector proportional to the data size to synchronize information among machines. Usually this step is the bottleneck of distributed machine learning. With fewer data passes and thus fewer rounds of communication, our method can be expected to have shorter training time than other algorithms in distributed environments.

8. Conclusions

In this paper, we propose a new optimization algorithm that utilizes the Hessian information to combine different update directions. Theoretically, our method has global linear and local quadratic convergence, while the cost per iteration is not expensive for solving ERM problems. Empirically, extensive experimental results show that for solving ERM problems, our method is competitive with state-of-the-art methods in terms of training time, but outperforms them in terms of data passes. An interesting future work is to consider possible choices of common directions other than the gradients. We also plan to exploit our method for distributed and disk-level training.

Acknowledgments

The implementation of Nesterov’s accelerated gradient is modified from an unreleased C# experimental code provided by Lin Xiao.

References

- Bernhard E. Boser, Isabelle Guyon, and Vladimir Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pages 144–152. ACM Press, 1992.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- Sébastien Bubeck, Yin Tat Lee, and Mohit Singh. A geometric alternative to Nesterov’s accelerated gradient descent. Technical report, 2015. arXiv preprint arXiv:1506.08187.
- Kai-Wei Chang and Dan Roth. Selective block minimization for faster convergence of limited memory large-scale linear models. In *Proceedings of the Seventeenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2011.
- John E. Dennis, Jr and Jorge J. Moré. Quasi-Newton methods, motivation and theory. *SIAM Review*, 19(1):46–89, 1977.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. LIBLINEAR: a library for large linear classification. *Journal of Machine Learning Research*, 9: 1871–1874, 2008. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/liblinear.pdf>.
- Reeves Fletcher and Colin M Reeves. Function minimization by conjugate gradients. *The Computer Journal*, 7(2):149–154, 1964.
- Euhanna Ghadimi, Hamid Reza Feyzmahdavian, and Mikael Johansson. Global convergence of the heavy-ball method for convex optimization. Technical report, 2014. arXiv preprint arXiv:1412.7457.
- William W Hager and Hongchao Zhang. A survey of nonlinear conjugate gradient methods. *Pacific Journal of Optimization*, 2(1):35–58, 2006.
- S. Sathiya Keerthi and Dennis DeCoste. A modified finite Newton method for fast solution of large scale linear SVMs. *Journal of Machine Learning Research*, 6:341–361, 2005.
- James E. Kelley, Jr. The cutting-plane method for solving convex programs. *Journal of the Society for Industrial and Applied Mathematics*, 1960.
- Laurent Lessard, Benjamin Recht, and Andrew Packard. Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1): 57–95, 2016.
- Chih-Jen Lin and Jorge J. Moré. Newton’s method for large-scale bound constrained problems. *SIAM Journal on Optimization*, 9:1100–1127, 1999.
- Chih-Jen Lin, Ruby C. Weng, and S. Sathiya Keerthi. Trust region Newton method for large-scale logistic regression. *Journal of Machine Learning Research*, 9:627–650, 2008. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/logistic.pdf>.
- Dong C. Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1):503–528, 1989.
- Olvi L. Mangasarian. A finite Newton method for classification. *Optimization Methods and Software*, 17(5):913–929, 2002.

Yurii E. Nesterov. A method of solving a convex programming problem with convergence rate $o(1/k^2)$. *Soviet Mathematics Doklady*, 27:372–376, 1983.

Yurii E. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Kluwer Academic Publishers, 2003.

Yurii E. Nesterov. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013.

Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer, second edition, 2006.

Boris Teodorovich Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.

Choon Hui Teo, S.V.N. Vishwanathan, Alex Smola, and Quoc V. Le. Bundle methods for regularized risk minimization. *Journal of Machine Learning Research*, 11:311–365, 2010.

Vladimir Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York, NY, 1995.

Hsiang-Fu Yu, Cho-Jui Hsieh, Kai-Wei Chang, and Chih-Jen Lin. Large linear classification when data cannot fit in memory. *ACM Transactions on Knowledge Discovery from Data*, 5(4):23:1–23:23, February 2012. URL http://www.csie.ntu.edu.tw/~cjlin/papers/kdd_disk_decomposition.pdf.

Appendix A. Analysis of Common-directions Method with Single Inner Iteration

We will show that under Assumptions 2.1, Algorithm 1 converges linearly by Theorem 2.3. If in addition Assumption 2.2 is satisfied, Algorithm 1 converges quadratically by Theorem 2.4.

Lemma A.1 *Let P be an n by m real matrix such that $P^\top P = I$. Then*

$$\|P^\top \mathbf{v}\| \leq \|\mathbf{v}\|, \quad \forall \mathbf{v} \in \mathbb{R}^n.$$

Further, equality holds if and only if

$$\mathbf{v} \in \text{span}(P) \equiv \{P\mathbf{t} \mid \forall \mathbf{t} \in \mathbb{R}^m\}.$$

Proof Note that \mathbf{v} can be decomposed as

$$\mathbf{v} = PP^\top \mathbf{v} + (I - PP^\top)\mathbf{v}. \tag{38}$$

Taking norm at both sides and applying $P^\top P = I$, there is

$$\|\mathbf{v}\|^2 = \|P^\top \mathbf{v}\|^2 + \|(I - PP^\top)\mathbf{v}\|^2 \geq \|P^\top \mathbf{v}\|^2. \tag{39}$$

Thus, the first result follows. For the second result, we know that \mathbf{v} is in the span of P if and only if its orthogonal projection to the span of P is itself. Therefore,

$$\mathbf{v} \in \text{span}(P) \iff \mathbf{v} = PP^\top \mathbf{v}.$$

Combining this with (38), we obtain

$$\mathbf{v} \in \text{span}(P) \implies (I - PP^\top)\mathbf{v} = \mathbf{0} \implies \|P^\top \mathbf{v}\| = \|\mathbf{v}\|.$$

Similarly, with (39), we have

$$\|P^\top \mathbf{v}\| = \|\mathbf{v}\| \implies \|(I - PP^\top)\mathbf{v}\| = 0 \implies \mathbf{v} \in \text{span}(P).$$

Thus, the second result follows from the above two statements. \blacksquare

Lemma A.2 *For every iteration k of Algorithm 1, we have*

$$P_k^\top P_k = I, \quad (40)$$

and

$$\|P_k^\top \nabla f(\mathbf{w}_k)\| = \|\nabla f(\mathbf{w}_k)\|. \quad (41)$$

Proof The augmenting of P in Algorithm 1 is simply the Gram-Schmidt process, and thus (40) holds. Further, from the fact that $(I - P_{k-1}P_{k-1}^\top)\nabla f(\mathbf{w}_k)$ is used to generate P_k , we have that $\nabla f(\mathbf{w}_k) \in \text{span}(P_k)$, so (41) follows from Lemma A.1. \blacksquare

Lemma A.3 *Under Assumption 2.1,*

$$\sigma I \preceq P_k^\top \nabla^2 f(\mathbf{w}_k) P_k \preceq \rho I. \quad (42)$$

Proof Assumption 2.1 indicates that

$$\sigma I \preceq \nabla^2 f(\mathbf{w}) \preceq \rho I, \quad \forall \mathbf{w}. \quad (43)$$

Combining (43) and (40) in Lemma A.2, one has that for any \mathbf{t} ,

$$\begin{aligned} \mathbf{t}^\top P_k^\top \nabla^2 f(\mathbf{w}_k) P_k \mathbf{t} &\geq \sigma \|P_k \mathbf{t}\|^2 = \sigma \|\mathbf{t}\|^2, \\ \mathbf{t}^\top P_k^\top \nabla^2 f(\mathbf{w}_k) P_k \mathbf{t} &\leq \rho \|P_k \mathbf{t}\|^2 = \rho \|\mathbf{t}\|^2, \end{aligned} \quad (44)$$

which then imply (42). \blacksquare

Lemma A.4 *Under Assumption 2.1,*

$$\frac{1}{2\rho} \|\nabla f(\mathbf{w})\|^2 \leq f(\mathbf{w}) - f^* \leq \frac{1}{2\sigma} \|\nabla f(\mathbf{w})\|^2, \quad (45)$$

where f^* is the optimal objective value.

Proof The first inequality is from Nesterov (2003, Theorem 2.1.5) together with that

$$\nabla f(\mathbf{w}^*) = \mathbf{0},$$

where \mathbf{w}^* is the optimal solution. The second inequality follows from Nesterov (2003, Theorem 2.1.10). \blacksquare

Proof of Theorem 2.3

We will show that the update direction does not deviate from the steepest descent direction too much, and the line search procedure will then ensure the sufficient decrease of the function value for linear convergence. Let

$$\mathbf{t} = -(P_k^\top \nabla^2 f(\mathbf{w}_k) P_k)^{-1} P_k^\top \nabla f(\mathbf{w}_k). \quad (46)$$

Note that $(P_k^\top \nabla^2 f(\mathbf{w}_k) P_k)$ is invertible because of Lemma A.3.

We first show that the step size produced in the line search procedure in Algorithm 2 is lower-bounded. From (43) via Assumption 2.1, for any $\bar{\theta} > 0$, we have

$$f(\mathbf{w}_k + \bar{\theta} P_k \mathbf{t}) \leq f(\mathbf{w}_k) + \bar{\theta} \nabla f(\mathbf{w}_k)^\top P_k \mathbf{t} + \frac{\rho}{2} \bar{\theta}^2 \|P_k \mathbf{t}\|^2. \quad (47)$$

If

$$-\bar{\theta} \nabla f(\mathbf{w}_k)^\top P_k \mathbf{t} - \frac{\rho}{2} \bar{\theta}^2 \|P_k \mathbf{t}\|^2 > \frac{\lambda}{2} \bar{\theta}^2 \|P_k \mathbf{t}\|^2, \quad (48)$$

then with (47) the stopping condition of Algorithm 2 is satisfied. For $\bar{\theta} > 0$, (48) is equivalent to

$$\bar{\theta} \leq \frac{-\nabla f(\mathbf{w}_k)^\top P_k \mathbf{t}}{\|P_k \mathbf{t}\|^2} \frac{2}{\rho + \lambda}. \quad (49)$$

From (40) of Lemma A.2,

$$\|P_k \mathbf{t}\| = \|\mathbf{t}\|. \quad (50)$$

Thus, by the definition of \mathbf{t} in (46) and Lemma A.3,

$$\frac{-\nabla f(\mathbf{w}_k)^\top P_k \mathbf{t}}{\|P_k \mathbf{t}\|^2} = \frac{\mathbf{t}^\top (P_k^\top \nabla^2 f(\mathbf{w}_k) P_k) \mathbf{t}}{\|\mathbf{t}\|^2} \geq \sigma, \quad (51)$$

whenever $\nabla f(\mathbf{w}_k) \neq \mathbf{0}$. Then we have that (48), or equivalently (49), holds whenever

$$\bar{\theta} \leq \frac{2\sigma}{\rho + \lambda}.$$

This ensures the finite termination of Algorithm 2 and guarantees that the generated step size θ_k satisfies

$$\theta_k \geq \min(1, \frac{2\beta\sigma}{\rho + \lambda}) \geq \frac{\beta\sigma}{\rho + \lambda}, \quad (52)$$

where the second inequality is from $\beta < 1$, $\lambda > 0$, and $\sigma \leq \rho$. Therefore, the line-search procedure terminates after at most $\lceil \log_\beta(\beta\sigma/(\rho + \lambda)) \rceil$ steps.

Next we prove the linear convergence. By (50), the definition of \mathbf{t} in (46), Lemma A.3, and (41) of Lemmas A.2, we have

$$\begin{aligned} \|P_k \mathbf{t}\|^2 &= \|\mathbf{t}\|^2 \\ &= \nabla f(\mathbf{w}_k)^\top P_k (P_k^\top \nabla^2 f(\mathbf{w}_k) P_k)^{-2} P_k^\top \nabla f(\mathbf{w}_k) \\ &\geq \frac{1}{\rho^2} \|P_k^\top \nabla f(\mathbf{w}_k)\| = \frac{1}{\rho^2} \|\nabla f(\mathbf{w}_k)\|^2. \end{aligned} \quad (53)$$

Therefore, by (52) and (53), we obtain from the line search stopping condition

$$f(\mathbf{w}_k) - f(\mathbf{w}_k + \theta_k P_k \mathbf{t}) \geq \frac{\lambda \beta^2 \sigma^2}{2(\rho + \lambda)^2} \|P_k \mathbf{t}\|^2 \geq \frac{\lambda \beta^2 \sigma^2}{2(\rho + \lambda)^2 \rho^2} \|\nabla f(\mathbf{w}_k)\|^2.$$

This inequality and Lemma A.4 then imply that

$$f(\mathbf{w}_{k+1}) - f^* \leq f(\mathbf{w}_k) - f^* - \frac{\lambda \beta^2 \sigma^2}{2(\rho + \lambda)^2 \rho^2} \|\nabla f(\mathbf{w}_k)\|^2 \leq \left(1 - \frac{\lambda \beta^2 \sigma^3}{(\rho + \lambda)^2 \rho^2}\right) (f(\mathbf{w}_k) - f^*).$$

This non-asymptotic linear convergence holds for any λ . Specifically, if we take $\lambda = \rho$, then we have the iteration complexity

$$O\left(\frac{\rho^3}{\sigma^3} \log\left(\frac{1}{\epsilon}\right)\right),$$

and the result is established. ■

Proof of Theorem 2.4

If $\nabla f(\mathbf{w}_0) = \mathbf{0}$, then the initial point is the optimal solution, and there is nothing to prove. Therefore, because $\nabla f(\mathbf{w}_0)/\|\nabla f(\mathbf{w}_0)\| \neq \mathbf{0}$ is included in P_0 , we start our discussion with $\text{rank}(P_k) \geq 1$. The procedure of adding columns to P_k in Algorithm 1 ensures that that $\text{rank}(P_k)$ is a monotonically increasing and bounded sequence with

$$1 \leq \text{rank}(P_k) \leq n, \quad \forall k \in \mathbb{N}.$$

Thus, there exists a k_0 such that

$$\text{rank}(P_k) = \text{rank}(P_{k_0}), \quad \forall k \geq k_0.$$

Together with $\mathbf{w}_{k+1} = \mathbf{w}_k + \theta_k P_k \mathbf{t}_k$ in Algorithm 1, this implies

$$\mathbf{w}_k - \mathbf{w}_0 \in \text{span}(P_{k_0}), \quad \forall k \geq 0.$$

From Theorem 2.3, we know that $f(\mathbf{w}_k)$ converges to $f(\mathbf{w}^*)$. This property and the strong convexity of f imply that \mathbf{w}_k also converges to \mathbf{w}^* globally because

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \geq \frac{\sigma}{2} \|\mathbf{w}_k - \mathbf{w}^*\|^2.$$

Therefore, we have $\mathbf{w}^* - \mathbf{w}_0 \in \text{span}(P_{k_0})$. This means that the minimization procedure takes update directions wholly in the space $\text{span}(P_{k_0})$, and each step in the algorithm is

equivalent to applying the Newton method with line search in the subspace. Consider the following equivalent reformulation

$$\mathbf{w}_k = \mathbf{w}_0 + P_{k_0} \hat{\mathbf{t}}_k, \mathbf{w}^* = \mathbf{w}_0 + P_{k_0} \hat{\mathbf{t}}^*, \forall k \geq k_0.$$

Our method after the k_0 -th iteration is equivalent to applying the Newton method over the variables $\hat{\mathbf{t}}$. That is, the following optimization problem is solved

$$\min_{\hat{\mathbf{t}}} f(\mathbf{w}_0 + P_{k_0} \hat{\mathbf{t}}). \quad (54)$$

From Lemma A.3, (54) possesses a unique optimal solution $\hat{\mathbf{t}}^*$. See also (46), which has the form of the Newton direction of (54). Note that for $g(\mathbf{t}) \equiv f(\mathbf{w}_0 + P_{k_0} \mathbf{t})$, Assumption 2.2 is still satisfied because with $\|P_{k_0}\| = 1$,

$$\begin{aligned} \|\nabla^2 g(\mathbf{t}_1) - \nabla^2 g(\mathbf{t}_2)\| &= \|P_{k_0}^\top (\nabla^2 f(\mathbf{w}_0 + P_{k_0} \mathbf{t}_1) - \nabla^2 f(\mathbf{w}_0 + P_{k_0} \mathbf{t}_2)) P_{k_0}\| \\ &\leq \|P_{k_0}\|^2 M \|P_{k_0}(\mathbf{t}_1 - \mathbf{t}_2)\| \leq M \|\mathbf{t}_1 - \mathbf{t}_2\|. \end{aligned}$$

Further, the line search process is scale-invariant because the columns in P_{k_0} are orthogonal due to the Schmidt process. Thus, the quadratic convergence of $\hat{\mathbf{t}}_k$ follows from that of the Newton method for problems satisfying Assumption 2.2 (Nocedal and Wright, 2006, Theorem 3.5).

$$\|\hat{\mathbf{t}}_{k+1} - \hat{\mathbf{t}}^*\| = O(\|\hat{\mathbf{t}}_k - \hat{\mathbf{t}}^*\|^2).$$

Because columns in P_{k_0} are orthogonal to each other, we have that

$$\|P_{k_0}(\hat{\mathbf{t}}_{k+1} - \hat{\mathbf{t}}^*)\| = \|(\hat{\mathbf{t}}_{k+1} - \hat{\mathbf{t}}^*)\| = O(\|\hat{\mathbf{t}}_k - \hat{\mathbf{t}}^*\|^2) = O(\|P_{k_0}(\hat{\mathbf{t}}_k - \hat{\mathbf{t}}^*)\|^2).$$

Therefore, quadratic convergence for the original variables \mathbf{w}_k holds. \blacksquare

Appendix B. Proof of Theorem 4.1

We first prove that, if every iteration of Algorithm 4 reaches the optimal solution of (21), then the optimal convergence rate is guaranteed. The main idea of this proof follows from the convergence rate proof of Nesterov's accelerated gradient (Nesterov, 1983) in Bubeck (2015, Theorem 3.17).

Let \mathbf{w}_k be the iterate obtained from Algorithm 4 at the k -th iteration. Define

$$\phi_{k+1}(\mathbf{w}) \equiv (1 - \alpha) \phi_k(\mathbf{w}) + \alpha \left(f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w} - \mathbf{w}_k) + \frac{\sigma}{2} \|\mathbf{w} - \mathbf{w}_k\|^2 \right), k \geq 0, \quad (55)$$

$$\text{with } \alpha \equiv \sqrt{\frac{\sigma}{\rho}} \in (0, 1], \quad \text{and} \quad \phi_0(\mathbf{w}) \equiv \frac{\sigma}{2} \|\mathbf{w} - \mathbf{w}_0\|^2 + f(\mathbf{w}_0). \quad (56)$$

We will first show that, by induction,

$$\nabla^2 \phi_k(\mathbf{w}) = \sigma I, \quad \forall k. \quad (57)$$

The claim holds trivially at $k = 0$ from (56). Now suppose the claim (57) holds for some $k \geq 0$. For the $(k + 1)$ -th iteration, by (55) we have

$$\nabla^2 \phi_{k+1}(\mathbf{w}) = (1 - \alpha) \nabla^2 \phi_k(\mathbf{w}) + \alpha \sigma I = \sigma I.$$

Thus, the induction on (57) is established.

Now we examine the recurrence relation in ϕ_k . Because $\phi_0(\mathbf{w})$ and the recurrence (55) are all strongly convex quadratic, we have $\phi_k(\mathbf{w})$ is also a strongly convex quadratic by induction. Let \mathbf{v}_k be a constant that minimizes ϕ_k , which implies $\nabla \phi_k(\mathbf{v}_k) = \mathbf{0}$. By (57) and the Taylor expansion, we obtain

$$\begin{aligned}\phi_k(\mathbf{w}) &= \phi_k(\mathbf{v}_k) + \nabla \phi_k(\mathbf{v}_k)^\top (\mathbf{w} - \mathbf{v}_k) + \frac{\sigma}{2} \|\mathbf{w} - \mathbf{v}_k\|^2 \\ &= \phi_k(\mathbf{v}_k) + \frac{\sigma}{2} \|\mathbf{w} - \mathbf{v}_k\|^2, \quad \forall \mathbf{w}.\end{aligned}\tag{58}$$

Further, let \mathbf{v}_{k+1} be the minimizer of ϕ_{k+1} . Using the derivative of (55) together with (58), we have

$$\begin{aligned}\nabla \phi_{k+1}(\mathbf{v}_{k+1}) &= (1 - \alpha) \nabla \phi_k(\mathbf{v}_{k+1}) + \alpha (\nabla f(\mathbf{w}_k) + \sigma(\mathbf{v}_{k+1} - \mathbf{w}_k)) \\ &= (1 - \alpha) \sigma(\mathbf{v}_{k+1} - \mathbf{v}_k) + \alpha (\nabla f(\mathbf{w}_k) + \sigma(\mathbf{v}_{k+1} - \mathbf{w}_k)) = \mathbf{0}.\end{aligned}$$

Thus, we obtain a recurrent definition of \mathbf{v}_{k+1} :

$$\begin{aligned}\mathbf{v}_0 &= \mathbf{w}_0, \\ \mathbf{v}_{k+1} &= (1 - \alpha) \mathbf{v}_k + \alpha \mathbf{w}_k - \frac{\alpha}{\sigma} \nabla f(\mathbf{w}_k).\end{aligned}\tag{59}$$

From (59), we have

$$\mathbf{v}_{k+1} - \mathbf{v}_k = -\alpha(\mathbf{v}_k - \mathbf{w}_k) - \frac{\alpha}{\sigma} \nabla f(\mathbf{w}_k),\tag{60}$$

$$\mathbf{v}_{k+1} - \mathbf{w}_k = (1 - \alpha)(\mathbf{v}_k - \mathbf{w}_k) - \frac{\alpha}{\sigma} \nabla f(\mathbf{w}_k).\tag{61}$$

Next, we obtain a recurrent relation of $\phi_{k+1}(\mathbf{v}_{k+1})$.

$$\begin{aligned}\phi_{k+1}(\mathbf{v}_{k+1}) &= (1 - \alpha) \left(\phi_k(\mathbf{v}_k) + \frac{\sigma}{2} \|\mathbf{v}_{k+1} - \mathbf{v}_k\|^2 \right) \\ &\quad + \alpha \left(f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{v}_{k+1} - \mathbf{w}_k) + \frac{\sigma}{2} \|\mathbf{v}_{k+1} - \mathbf{w}_k\|^2 \right)\end{aligned}\tag{62}$$

$$\begin{aligned}&= (1 - \alpha) \phi_k(\mathbf{v}_k) + \alpha f(\mathbf{w}_k) \\ &\quad + \alpha(1 - \alpha) \left(\nabla f(\mathbf{w}_k)^\top (\mathbf{v}_k - \mathbf{w}_k) + \frac{\sigma}{2} \|\mathbf{v}_k - \mathbf{w}_k\|^2 \right) - \frac{1}{2\rho} \|\nabla f(\mathbf{w}_k)\|^2,\end{aligned}\tag{63}$$

where (62) is by (55) and (58), and (63) is by (60), (61), and applying the definition $\alpha = \sqrt{\sigma/\rho}$ to the coefficient of $\|\nabla f(\mathbf{w}_k)\|^2$. Note that $\{\phi_k\}$ is called the estimation sequence in Nesterov (1983).

Now we show by induction that, for \mathbf{w}_k minimizing (21) exactly,

$$\min_{\mathbf{w} \in \mathbb{R}^n} \phi_k(\mathbf{w}) = \phi_k(\mathbf{v}_k) \geq f(\mathbf{w}_k), \quad \forall k \geq 0. \quad (64)$$

First, (64) holds trivially at $k = 0$ from the definition (56). Assume (64) holds at the k -th iteration. At the $(k + 1)$ -th iteration, (63) and (64) imply

$$\phi_{k+1}(\mathbf{v}_{k+1}) \geq f(\mathbf{w}_k) + \alpha(1 - \alpha) \left(\nabla f(\mathbf{w}_k)^\top (\mathbf{v}_k - \mathbf{w}_k) + \frac{\sigma}{2} \|\mathbf{v}_k - \mathbf{w}_k\|^2 \right) - \frac{1}{2\rho} \|\nabla f(\mathbf{w}_k)\|^2. \quad (65)$$

By the definition of P_k and induction on (59), both \mathbf{v}_k and \mathbf{w}_k belong to $\mathbf{w}_0 + \text{span}(P_{k-1})$. Thus, the optimality of \mathbf{w}_k on

$$\min_t f(\mathbf{w}_{k-1} + P_{k-1}\mathbf{t})$$

implies that

$$P_{k-1}^\top \nabla f(\mathbf{w}_k) = \mathbf{0}. \quad (66)$$

Thus

$$\nabla f(\mathbf{w}_k)^\top (\mathbf{v}_k - \mathbf{w}_k) = 0. \quad (67)$$

This result, (65), and $\alpha \in (0, 1]$ from (56) then lead to

$$\phi_{k+1}(\mathbf{v}_{k+1}) \geq f(\mathbf{w}_k) - \frac{1}{2\rho} \|\nabla f(\mathbf{w}_k)\|^2. \quad (68)$$

Further,

$$\begin{aligned} f(\mathbf{w}_k) - \frac{1}{2\rho} \|\nabla f(\mathbf{w}_k)\|^2 &= f(\mathbf{w}_k) - \frac{1}{\rho} \|\nabla f(\mathbf{w}_k)\|^2 + \frac{\rho}{2} \frac{1}{\rho^2} \|\nabla f(\mathbf{w}_k)\|^2 \\ &\geq f(\mathbf{w}_k - \frac{1}{\rho} \nabla f(\mathbf{w}_k)) \end{aligned} \quad (69)$$

$$\geq f(\mathbf{w}_{k+1}). \quad (70)$$

The inequality (69) comes from the Taylor expansion of $f(\mathbf{w}_k - \nabla f(\mathbf{w}_k)/\rho)$ and (43):

$$\begin{aligned} f(\mathbf{w}_k - \frac{1}{\rho} \nabla f(\mathbf{w}_k)) &= f(\mathbf{w}_k) - \frac{1}{\rho} \|\nabla f(\mathbf{w}_k)\|^2 + \frac{1}{2\rho^2} \nabla f(\mathbf{w}_k)^\top \nabla^2 f(\tilde{\mathbf{w}}) \nabla f(\mathbf{w}_k) \\ &\leq f(\mathbf{w}_k) - \frac{1}{\rho} \|\nabla f(\mathbf{w}_k)\|^2 + \frac{1}{2\rho} \|\nabla f(\mathbf{w}_k)\|^2, \end{aligned}$$

where $\tilde{\mathbf{w}}$ is between \mathbf{w}_k and $\mathbf{w}_k - \frac{1}{\rho} \nabla f(\mathbf{w}_k)$. For (70), it holds because $\nabla f(\mathbf{w}_k) \in \text{span}(P_k)$ and \mathbf{w}_{k+1} is a minimizer of f over $\mathbf{w}_0 + \text{span}(P_k)$. Combining (68) and (70), we have (64) holds at iteration $k + 1$. Therefore, (64) holds for every $k \geq 0$ by induction.

Now, we arrive at the final steps. From (64), (55), and the strong convexity of f in Assumption 2.1, we have

$$\begin{aligned} & f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \\ & \leq \min_{\mathbf{w}} \phi_{k+1}(\mathbf{w}) - f(\mathbf{w}^*) \\ & \leq \phi_{k+1}(\mathbf{w}^*) - f(\mathbf{w}^*) \end{aligned} \tag{71}$$

$$\begin{aligned} & = (1-\alpha)\phi_k(\mathbf{w}^*) + \alpha(f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}^* - \mathbf{w}_k) + \frac{\sigma}{2}\|\mathbf{w}^* - \mathbf{w}_k\|^2) - f(\mathbf{w}^*) \\ & \leq (1-\alpha)\phi_k(\mathbf{w}^*) + \alpha f(\mathbf{w}^*) - f(\mathbf{w}^*) \\ & = (1-\alpha)(\phi_k(\mathbf{w}^*) - f(\mathbf{w}^*)) \\ & \leq (1-\alpha)^{k+1}(\phi_0(\mathbf{w}^*) - f(\mathbf{w}^*)). \end{aligned} \tag{72}$$

$$\tag{73}$$

Note that (73) comes from repeating the steps between (71) and (72). By taking logarithm at both sides, we have proven that Algorithm 4 achieves the optimal rate when (21) is solved exactly.

Next we discuss the situation when (21) is approximately solved. In the above proof we can see that having a \mathbf{w}_{k+1} that satisfies (64) is sufficient for proving the linear convergence. Because (64) comes from (65)-(68) and (70),¹⁰ we can consider running inner iterations in Algorithm 4 until the following inner stopping condition is satisfied.

$$\left\{ \begin{array}{l} \nabla f(\mathbf{w}_{k+1})^\top (\mathbf{v}_{k+1} - \mathbf{w}_{k+1}) + \frac{\sigma}{2}\|\mathbf{v}_{k+1} - \mathbf{w}_{k+1}\|^2 \geq 0, \end{array} \right. \tag{74a}$$

$$\left\{ \begin{array}{l} f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k - \frac{1}{\rho}\nabla f(\mathbf{w}_k)). \end{array} \right. \tag{74b}$$

Denote \mathbf{t}_k^* as the exact solution of (21) and let

$$\mathbf{w}_{k+1}^* \equiv \mathbf{w}_k + P_k \mathbf{t}_k^*.$$

From the discussion in (67) and (70), clearly \mathbf{w}_{k+1}^* satisfies (74). The remaining task is to prove that (74) can be achieved in a finite number of iterations. If \mathbf{w}_{k+1}^* strictly satisfies inequalities in (74), then by the continuity of f and ∇f , there is a neighborhood of \mathbf{w}_{k+1}^* in which every point satisfies (74). Because inner iterations lead to the convergence to \mathbf{w}_{k+1}^* , the inner loop must terminate finitely. On the other hand, we can easily rule out situations where \mathbf{w}_{k+1}^* satisfies one of the conditions in (74) as an equality. When (74a) becomes an equality under \mathbf{w}_{k+1}^* , (67) implies that

$$\mathbf{w}_{k+1}^* = \mathbf{v}_{k+1}.$$

When (74b) becomes an equality, the strong convexity of $f(\mathbf{w}_k + P_k \mathbf{t})$ on \mathbf{t} , which follows from Lemma A.3, implies that the optimum \mathbf{w}_{k+1}^* of $\min_{\mathbf{t}} f(\mathbf{w}_k + P_k \mathbf{t})$ is unique. This property, the fact that $\nabla f(\mathbf{w}_k)$ lies in the column space of P_k , and the equality in (74b) imply that

$$\mathbf{w}_{k+1}^* = \mathbf{w}_k - \frac{1}{\rho}\nabla f(\mathbf{w}_k).$$

10. Note that we ensure (74) holds for every k , so (74a) with k rather than $k+1$ is what we use for (65).

Therefore, before starting the inner loop of Algorithm 4, we can first check if \mathbf{v}_{k+1} or $\mathbf{w}_k - (1/\rho)\nabla f(\mathbf{w}_k)$ satisfies (74). If one of them satisfies both inequalities, then we can take this vector as \mathbf{w}_{k+1} without going into the loop. Otherwise, neither \mathbf{v}_{k+1} nor $\mathbf{w}_k - (1/\rho)\nabla f(\mathbf{w}_k)$ is \mathbf{w}_{k+1}^* , so \mathbf{w}_{k+1}^* must strictly satisfy inequalities in (74).

Therefore, after a finite number of inner iterations we can always find an iterate \mathbf{w}_{k+1} satisfying (74). This condition ensures (68) and (70) hold. Therefore, (64) is established.

Appendix C. Proof of Theorem 4.2

At the k -th iteration, if (21) is solved exactly, the optimality condition (66) implies that any column in P_k is orthogonal to $\nabla f(\mathbf{w}_{k+1})$. Therefore, if $\nabla f(\mathbf{w}_{k+1}) \neq \mathbf{0}$, the new gradient will augment P_{k+1} by one column. By induction, if $\nabla f(\mathbf{w}_{k'}) \neq \mathbf{0}$, $\forall k' = 0, \dots, k-1$, then P_k has k orthogonal columns. Because each column of P_k is in \mathbb{R}^n , there can be at most n orthogonal directions in P_k . Therefore, when $k \geq n$, $\text{span}(P_k) = R^n$ and the solution of (21) is equivalent to the solution of (1). Thus, the optimum is reached within n iterations.

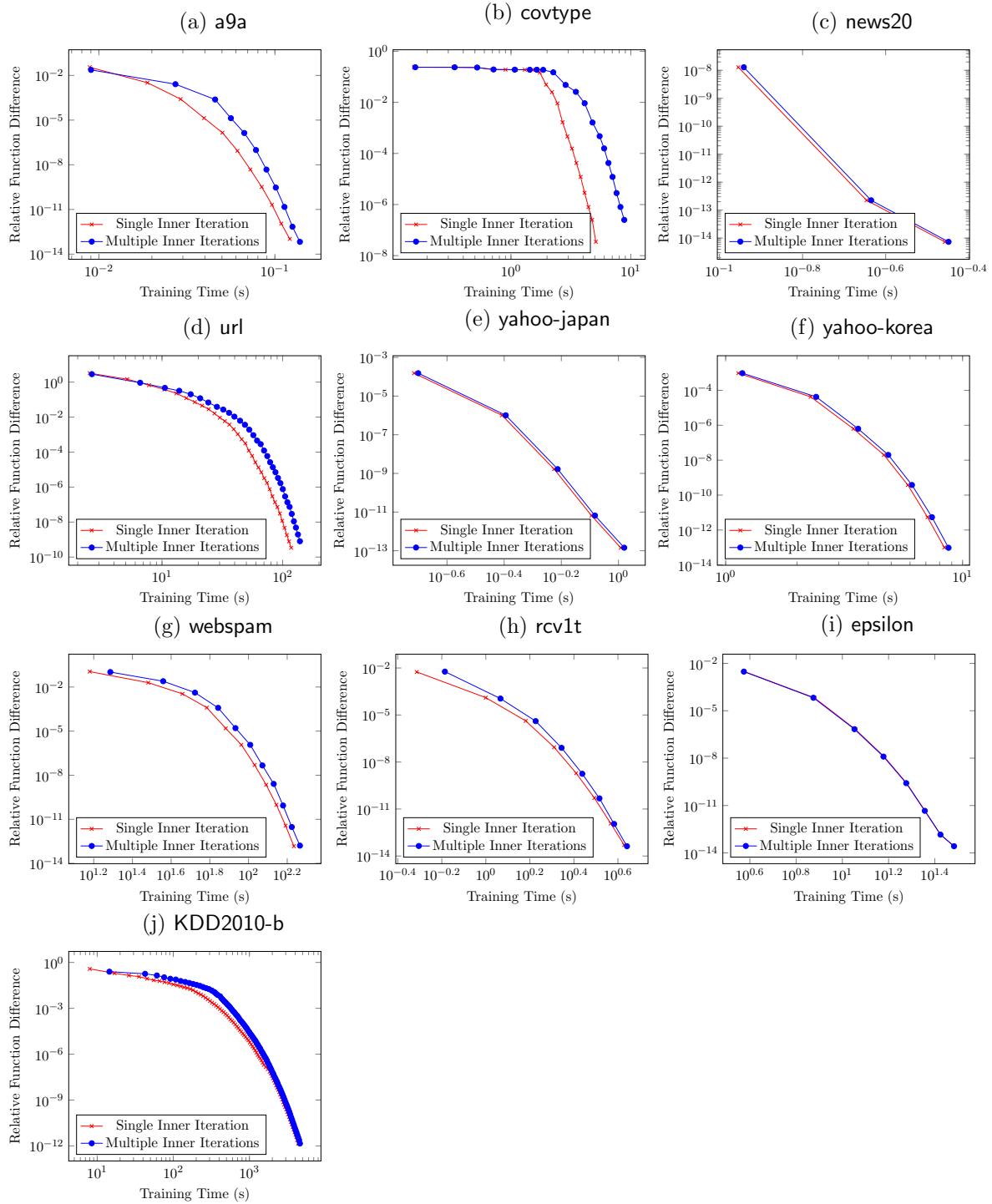


Figure 1: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present training time (in log scale) of logistic regression with $C = 10^{-3}$.

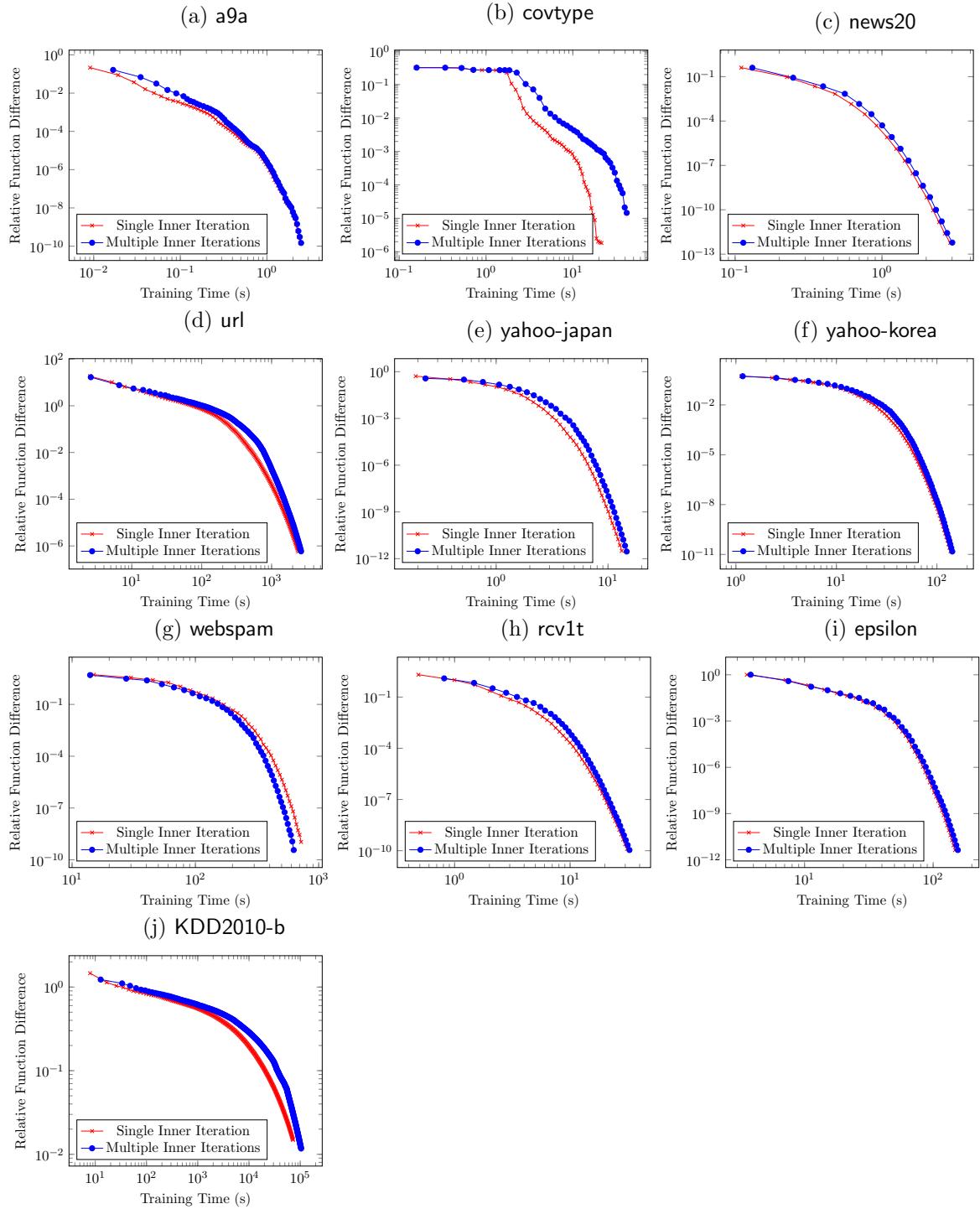


Figure 2: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present training time (in log scale) of logistic regression with $C = 10^1$.

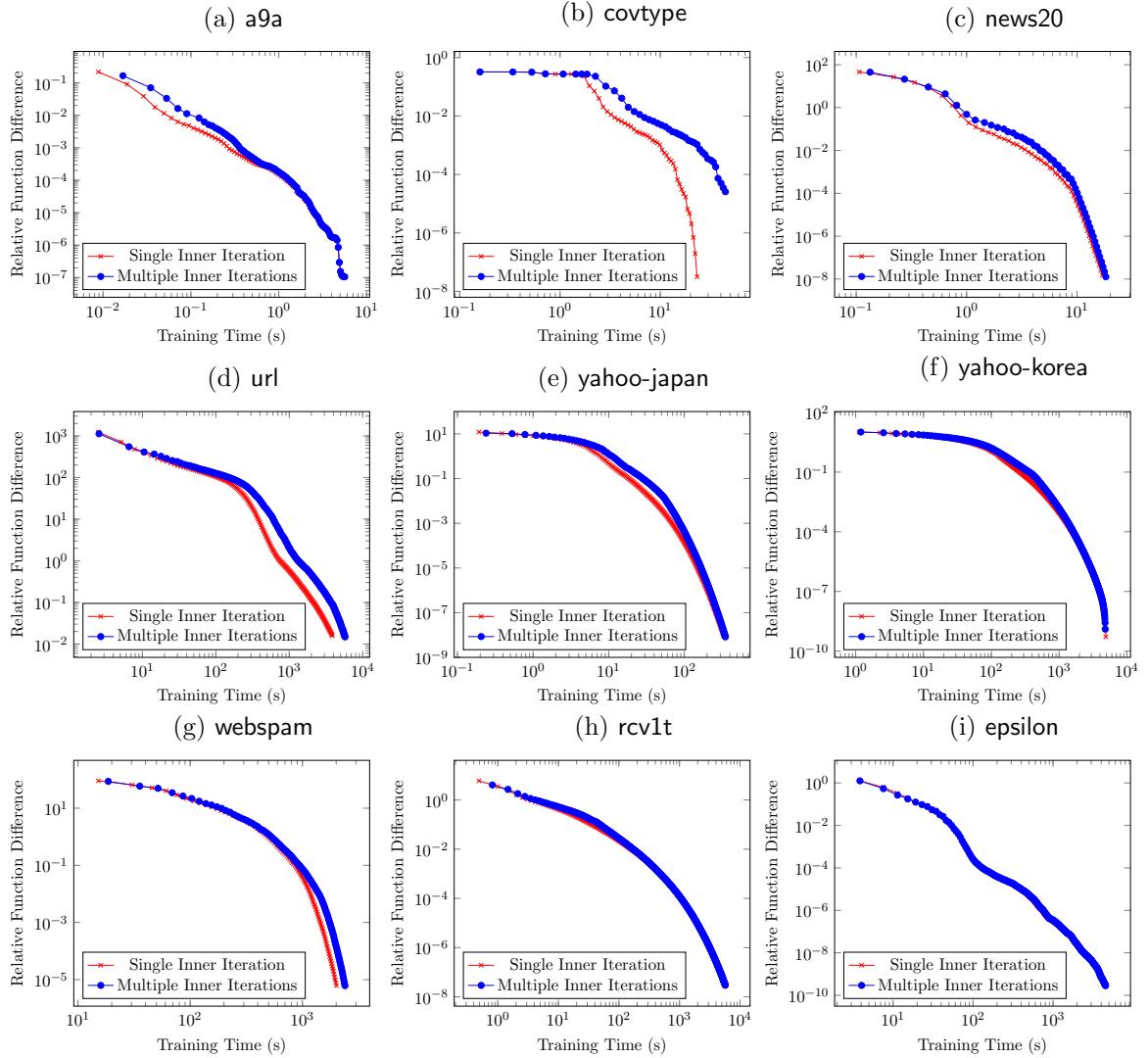


Figure 3: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present training time (in log scale) of logistic regression with $C = 10^3$.

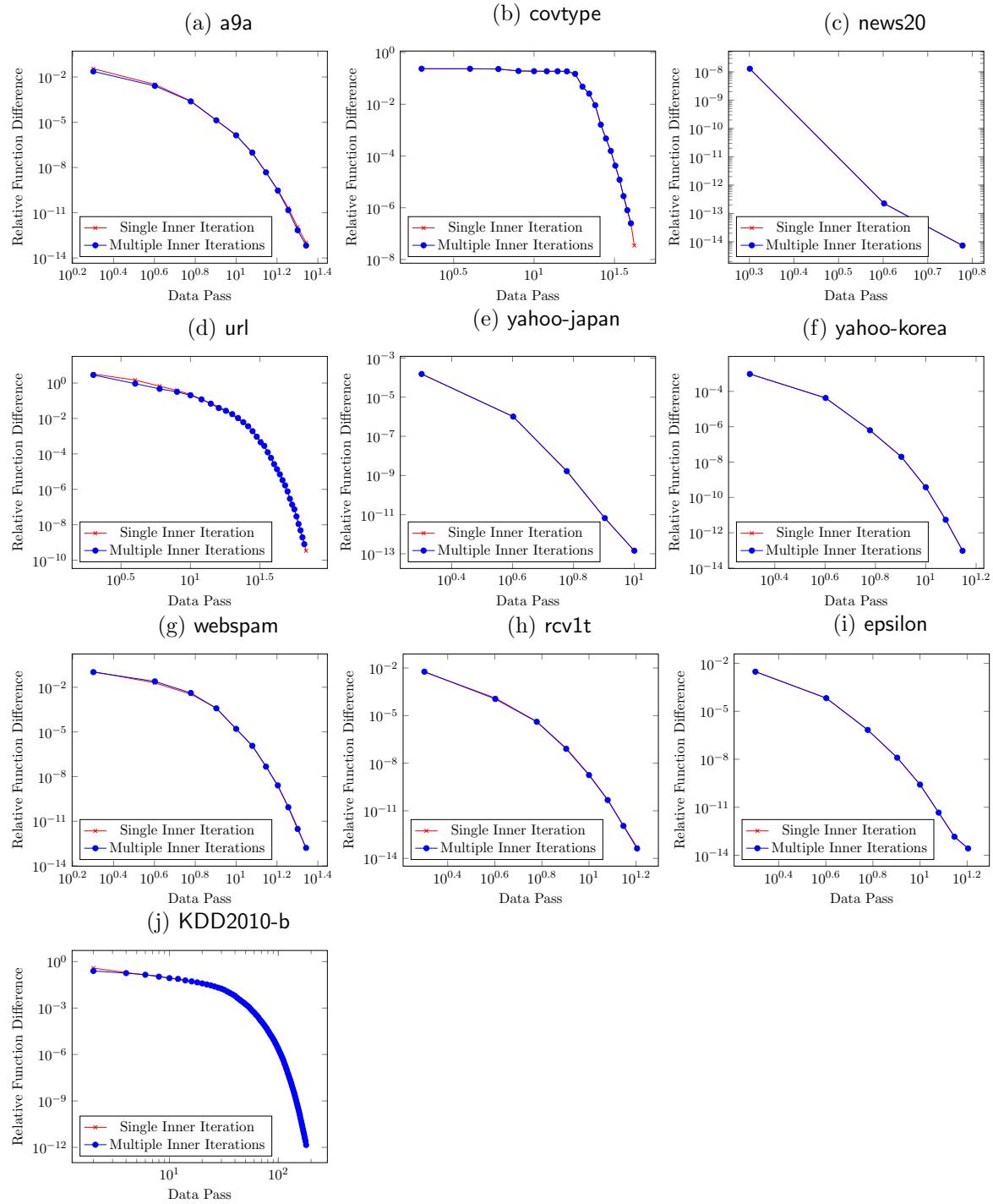


Figure 4: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present data passes (in log scale) of logistic regression with $C = 10^{-3}$.

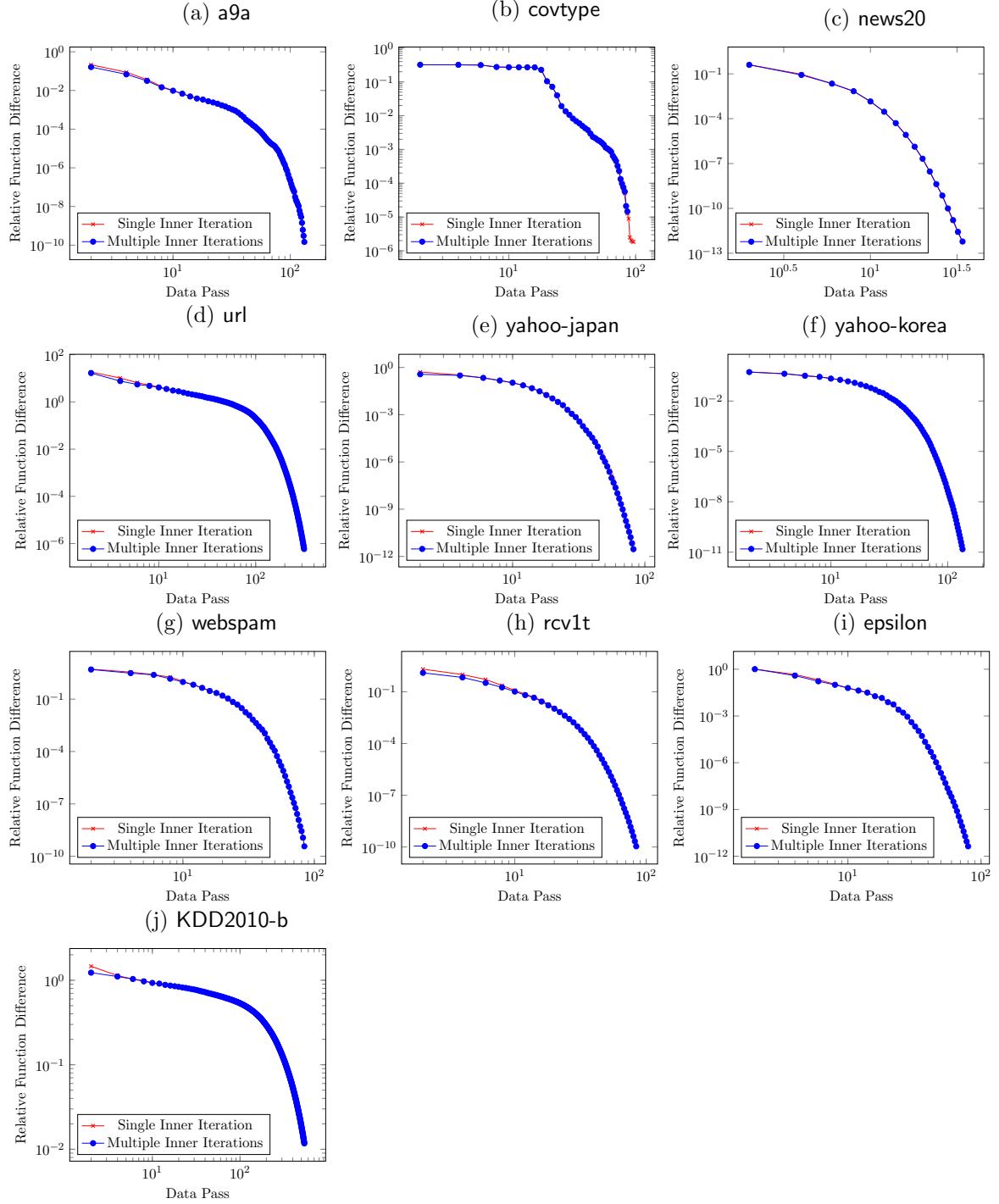


Figure 5: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present data passes (in log scale) of logistic regression with $C = 10^1$.

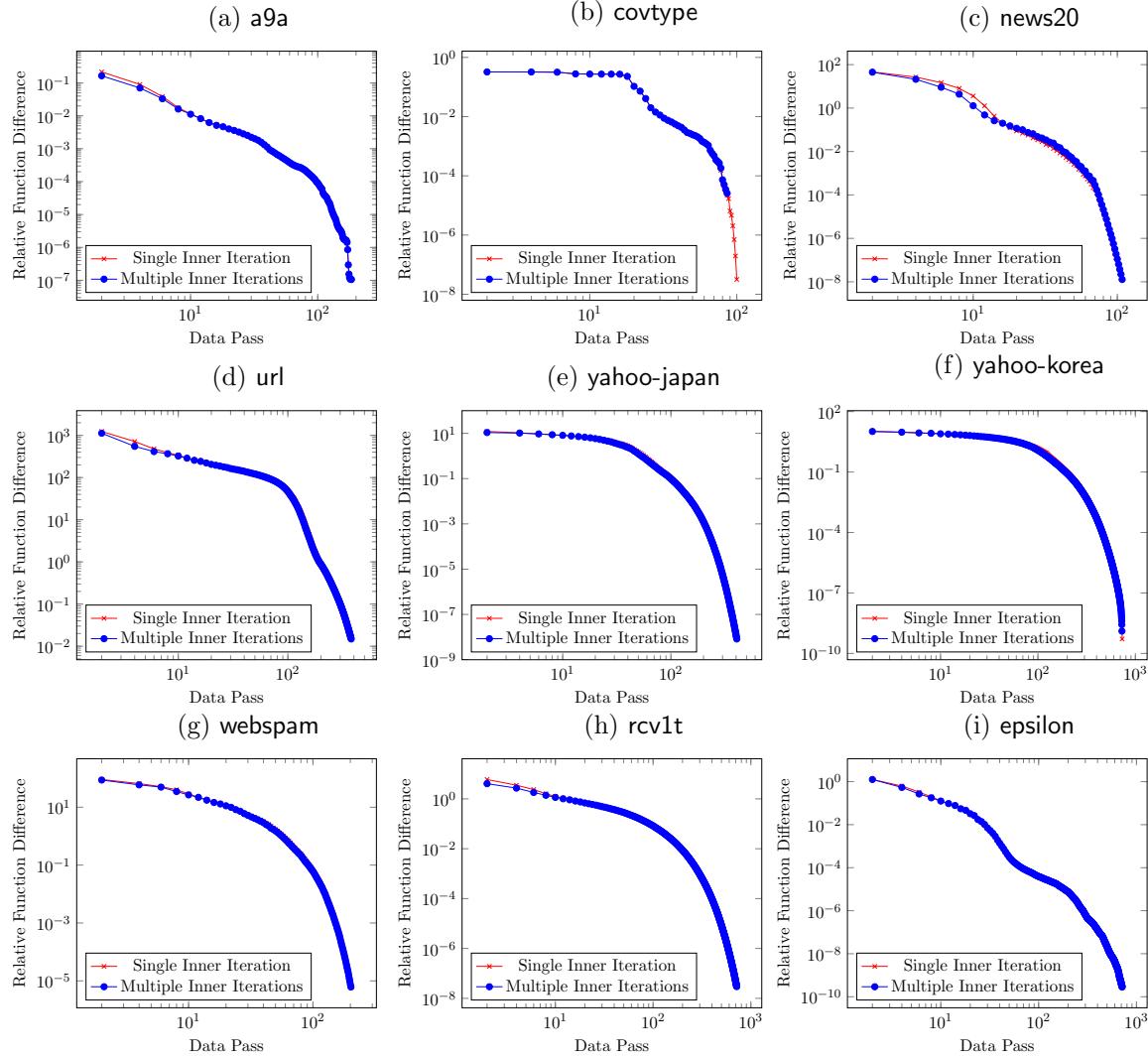


Figure 6: Comparison between single inner iteration and multiple inner iterations variants of the common-directions method. We present data passes (in log scale) of logistic regression with $C = 10^3$.

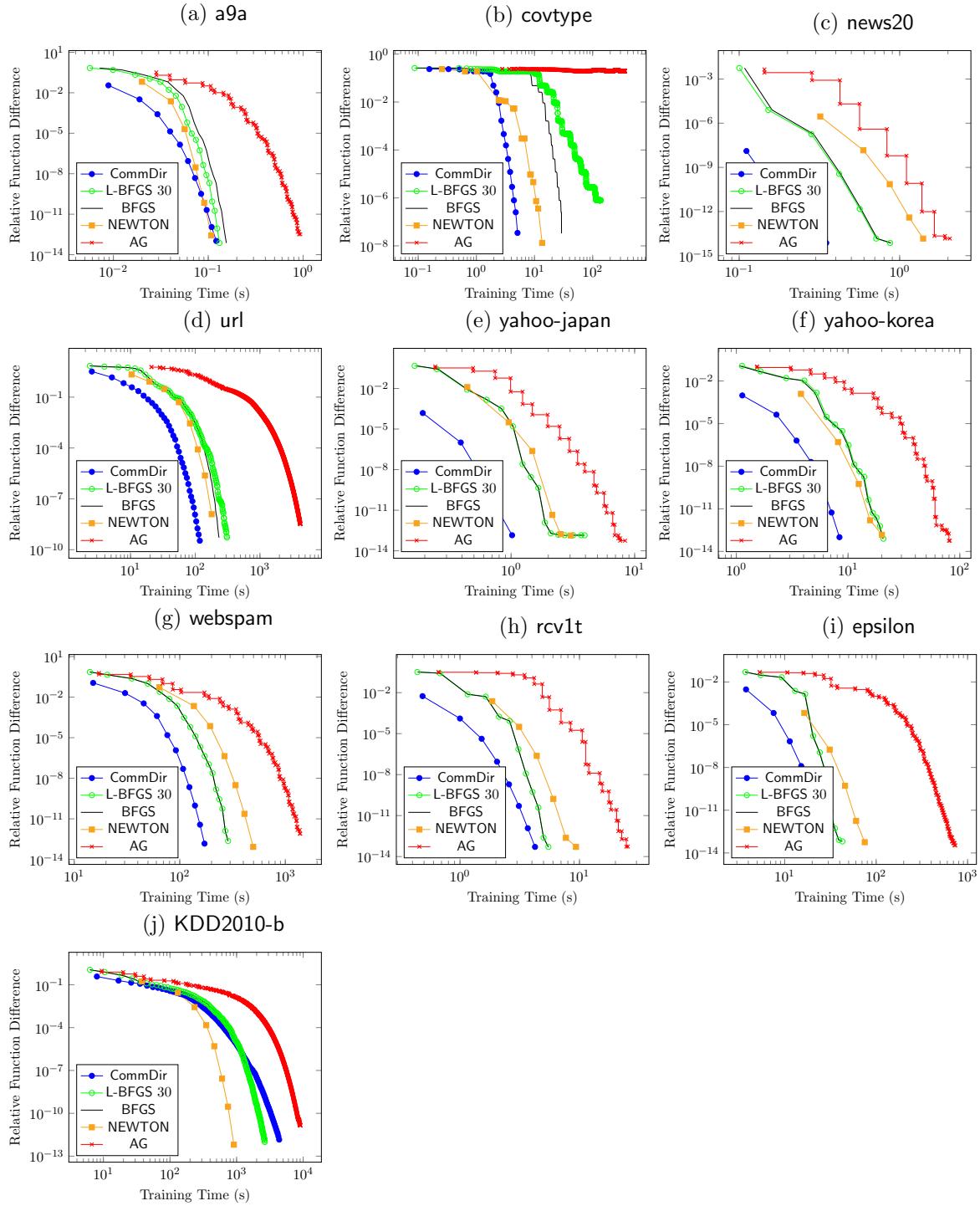


Figure 7: Training time of logistic regression with $C = 10^{-3}$.

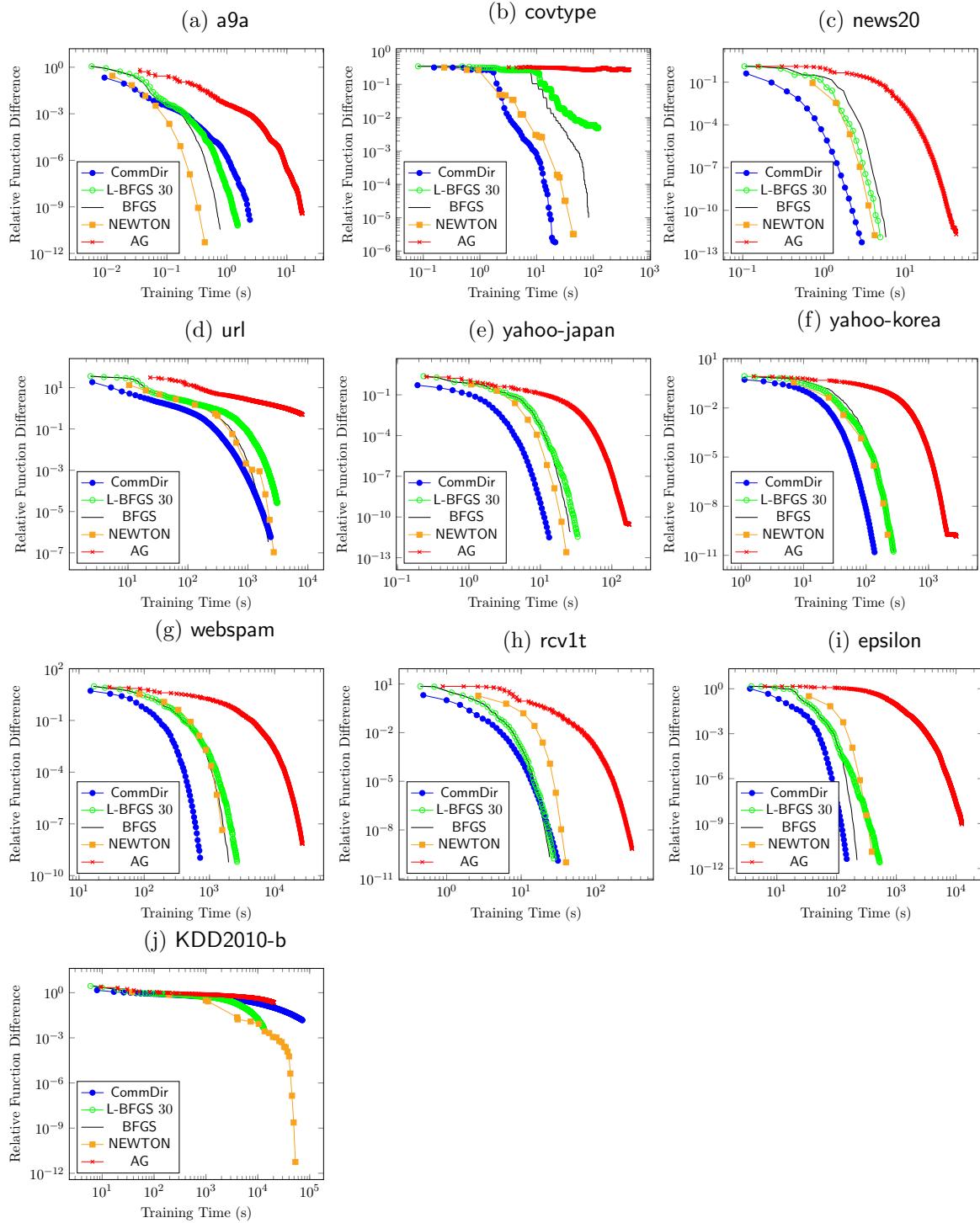


Figure 8: Training time of logistic regression with $C = 1$.

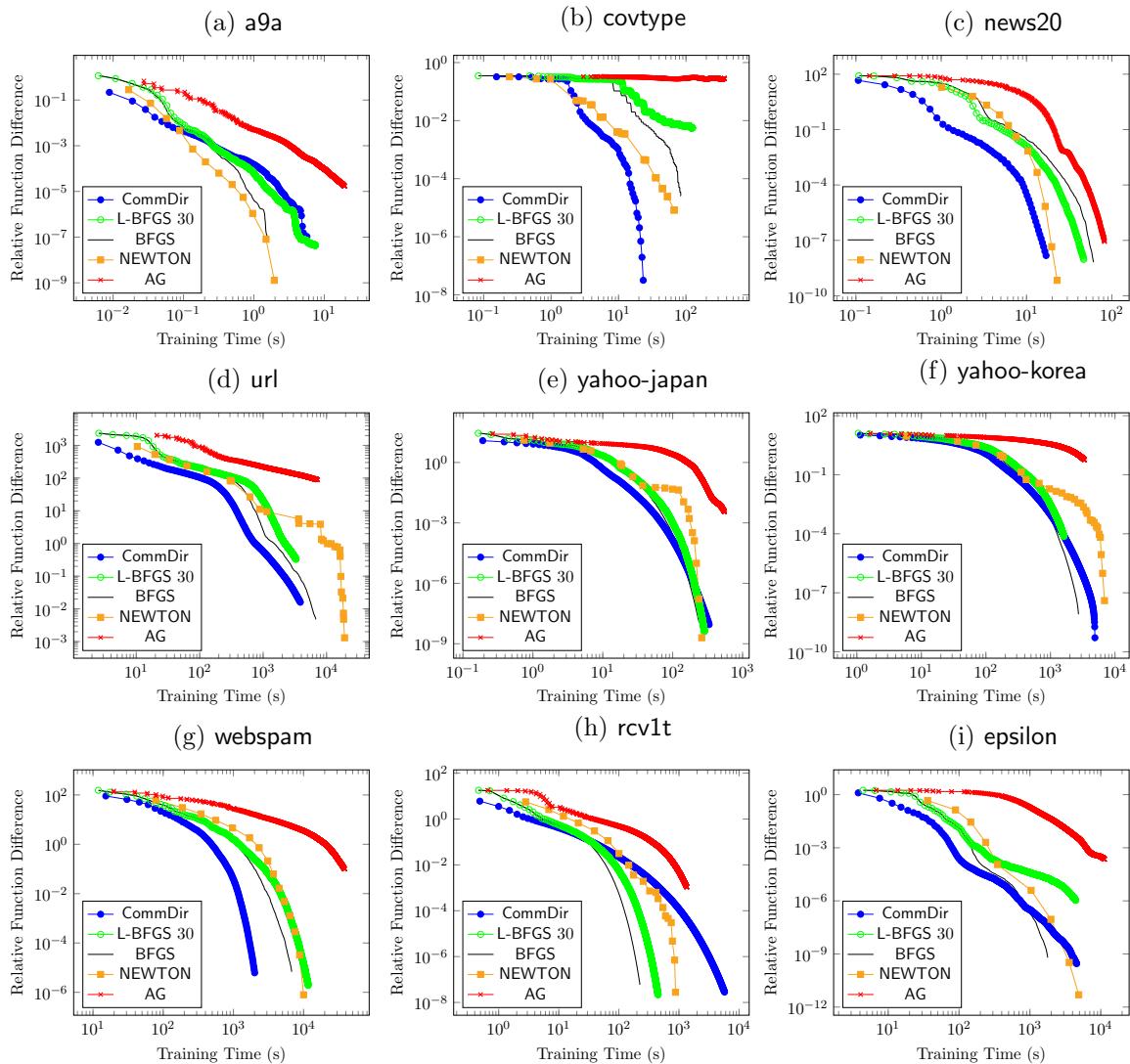


Figure 9: Training time of logistic regression with $C = 10^3$.

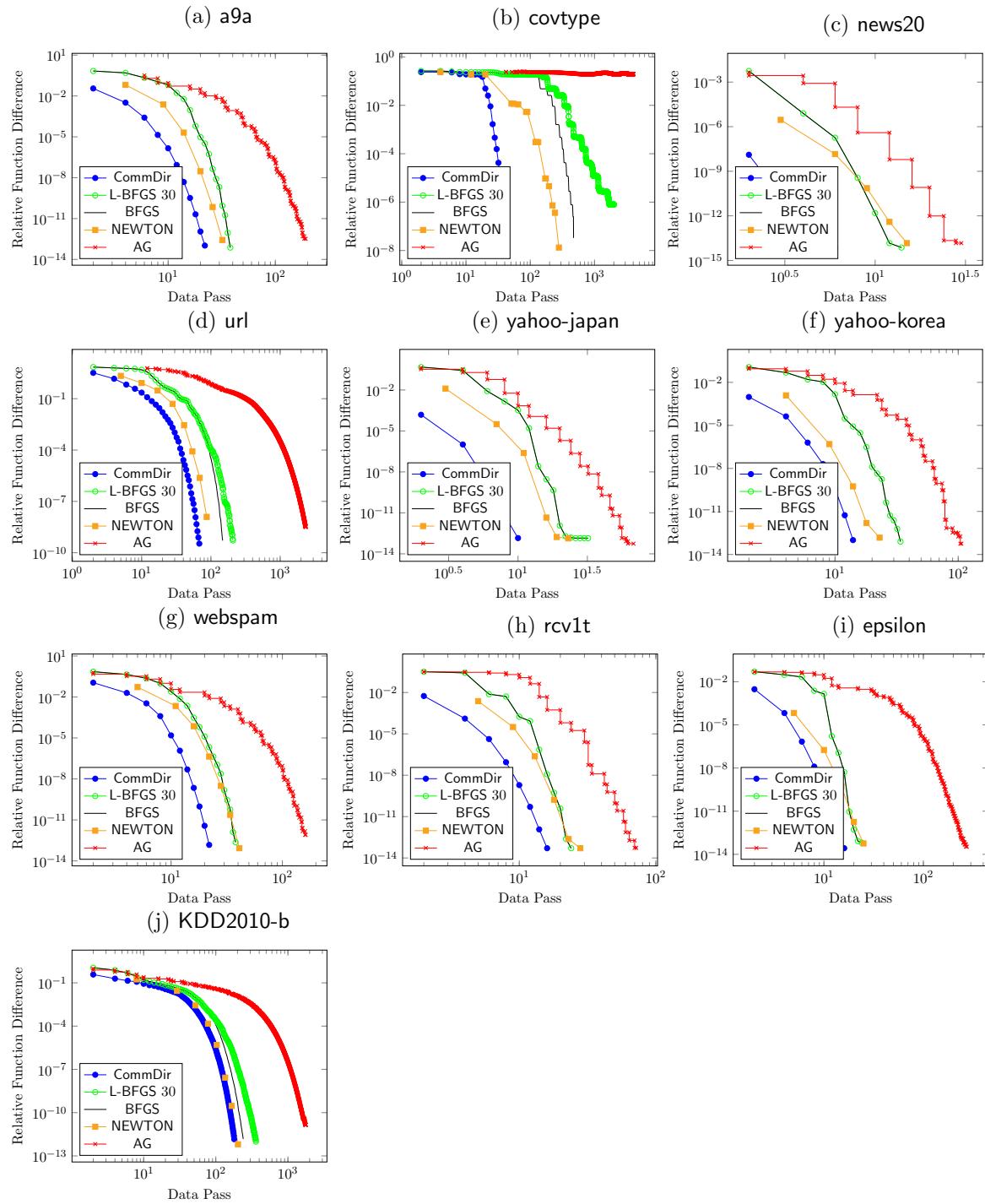


Figure 10: Number of data passes of logistic regression with $C = 10^{-3}$.

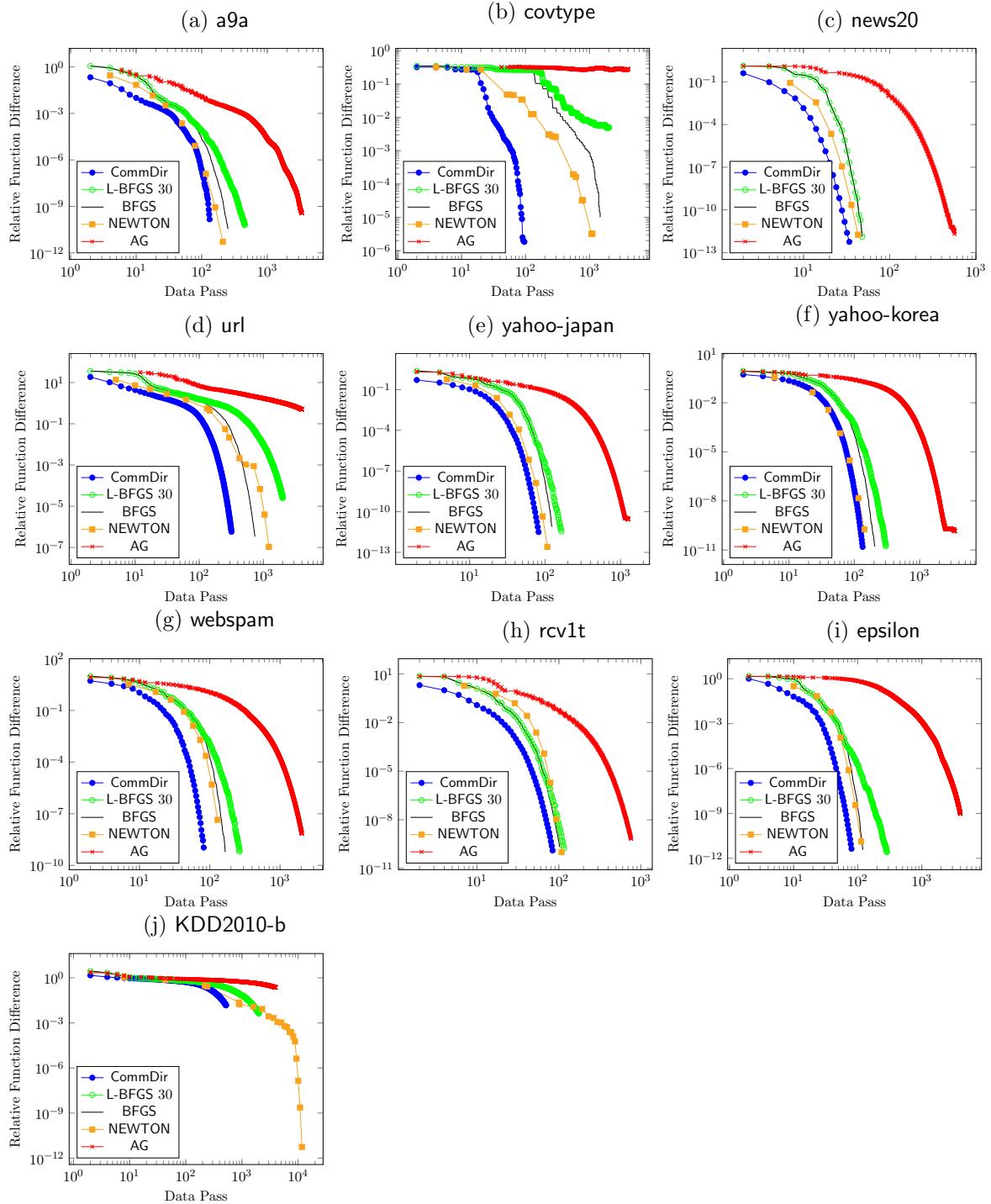


Figure 11: Number of data passes of logistic regression with $C = 1$.

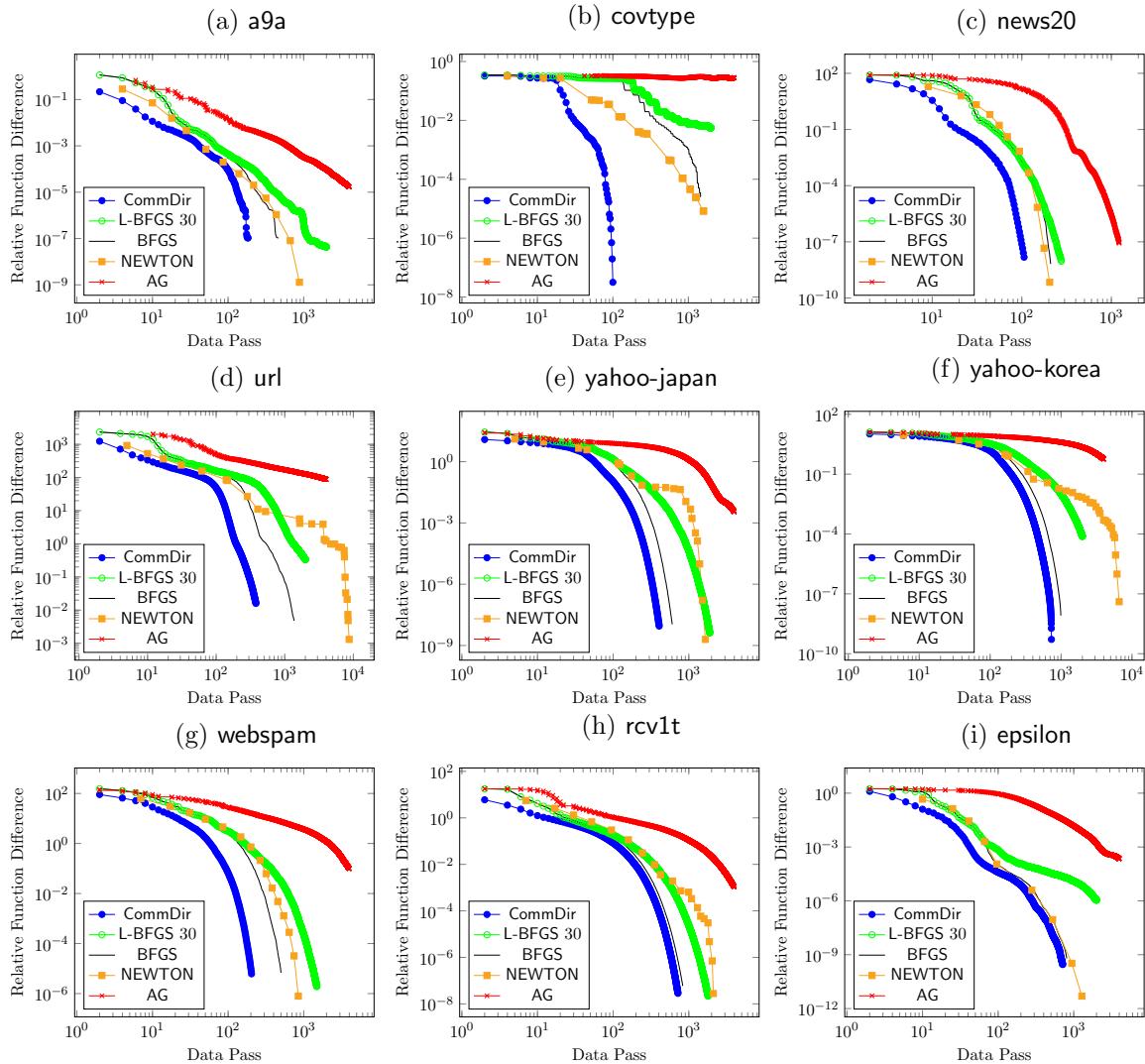


Figure 12: Number of data passes of logistic regression with $C = 10^3$.