

Preconditioned Conjugate Gradient Methods in Truncated Newton Frameworks for Large-scale Linear Classification

Chih-Yang Hsia

*Department of Computer Science
National Taiwan University*

R04922021@NTU.EDU.TW

Wei-Lin Chiang

*Department of Computer Science
National Taiwan University*

R06922166@CSIE.NTU.EDU.TW

Chih-Jen Lin

*Department of Computer Science
National Taiwan University*

CJLIN@CSIE.NTU.EDU.TW

Editors: Jun Zhu and Ichiro Takeuchi

Abstract

Truncated Newton method is one of the most effective optimization methods for large-scale linear classification. The main computational task at each Newton iteration is to approximately solve a quadratic sub-problem by an iterative procedure such as the conjugate gradient (CG) method. It is known that CG has slow convergence if the sub-problem is ill-conditioned. Preconditioned CG (PCG) methods have been used to improve the convergence of the CG method, but it is difficult to find a preconditioner that performs well in most situations. Further, because Hessian-free optimization techniques are incorporated for handling large data, many existing preconditioners are not directly applicable. In this work, we detailedly study some preconditioners that have been considered in past works for linear classification. We show that these preconditioners may not help to improve the training speed in some cases. After some investigation, we propose simple and effective techniques to make the PCG method more robust in a truncated Newton framework. The idea is to avoid the situation when a preconditioner leads to a much worse condition number than when it is not applied. We provide theoretical justification. Through carefully designed experiments, we demonstrate that our method can effectively reduce the training time for large-scale problems.

Keywords: large-scale linear classification, preconditioned conjugate gradient, Newton method

1. Introduction

In linear classification, logistic regression and linear SVM are two commonly used models. The model parameters, denoted as $\mathbf{w} \in \mathbb{R}^n$, are obtained by solving an unconstrained optimization problem

$$\min_{\mathbf{w}} f(\mathbf{w}).$$

Truncated Newton method is one of the most effective optimization methods for large-scale linear classification. The core computational task at the k th Newton iteration is to

approximately solve a sub-problem that is related to the following linear system

$$\nabla^2 f(\mathbf{w}_k) \mathbf{s}_k = -\nabla f(\mathbf{w}_k), \quad (1)$$

where $\nabla f(\mathbf{w}_k)$ and $\nabla^2 f(\mathbf{w}_k)$ are gradient and Hessian, respectively. For large-scale problems, the Hessian matrix is too large to be stored. Past works such as [Keerthi and DeCoste \(2005\)](#); [Lin et al. \(2008\)](#) have addressed this difficulty by showing that the special structure in linear classification allows us to solve (1) by a conjugate gradient (CG) procedure without forming the whole Hessian matrix. However, even with such a Hessian-free approach, the training of large-scale data sets is still a time-consuming process because of a possibly large number of CG steps.

It is well known that for ill-conditioned matrices, the convergence of CG methods may be slow. To reduce the number of CG steps, a well-known technique is the preconditioned conjugate gradient (PCG) method (e.g., [Concus et al., 1976](#); [Nash, 1985](#)). The idea is to apply a preconditioner matrix to possibly improve the condition of the linear system. Unfortunately, designing a suitable preconditioner is never an easy task. Because we need extra efforts to find and use the preconditioners, the cost per CG step becomes higher. Therefore, a smaller number of CG steps may not lead to shorter running time. Further, a PCG method is not guaranteed to reduce the number of CG steps.

For linear classification, some past works such as [Lin et al. \(2008\)](#); [Zhang and Lin \(2015\)](#); [Ma and Takáč \(2016\)](#) have applied PCG in the truncated Newton framework. However, whether PCG is useful remains a question to be studied. One challenge is that training a linear classifier is different from solving a single linear system in the following aspects.

1. Because a sequence of CG procedures is conducted, reducing the number of CG steps at one Newton iteration may not lead to the overall improvement. For example, suppose at the first iteration PCG takes fewer steps than CG but goes to a bad point for the overall optimization process. Then the total number of CG steps by using PCG may not be less.
2. The effectiveness of PCG may depend on when the optimization procedure is terminated. The work by [Lin et al. \(2008\)](#) concludes that a diagonal preconditioner is not consistently better than the ordinary CG. Later [Chin et al. \(2016\)](#) point out the conclusion in [Lin et al. \(2008\)](#) is based on a strict stopping condition in solving the optimization problem. If the optimization procedure is terminated earlier, a situation suitable for machine learning applications, PCG is generally useful. Therefore, careful designs are needed in evaluating the effectiveness of PCG.
3. Because the Hessian matrix is too large to be stored and we consider a Hessian-free approach, most existing preconditioners can not be directly used. This fact makes the application of PCG for linear classification very difficult.

In this work we show that existing preconditioners may not help to improve the training speed in some cases. After some investigation, we propose a simple and effective technique to make the PCG method more robust in a truncated Newton framework. The idea is to avoid the situation when a preconditioner leads to a much worse condition number than when it is not applied. The proposed technique can be applied in single-core, multi-core, or distributed settings because the running time is always proportional to the total number of CG steps.

This work is organized as follows. In Section 2, we introduce Newton methods for large-scale linear classification and discuss how PCG methods are applied to a trust-region Newton framework. In Section 3, we discuss several existing preconditioners and their possible weakness. We then propose a strategy in Section 4 to have a more robust preconditioner in the truncated Newton framework. Theoretical properties are investigated. In Section 5, we conduct thorough experiments and comparisons to show the effectiveness of our proposed setting. Proofs and additional experimental results are available in supplementary materials. The proposed method has been incorporated into the software LIBLINEAR (version 2.20 and after). Supplementary materials and experimental codes are at https://www.csie.ntu.edu.tw/~cjlin/papers/tron_pcg/.

2. Conjugate Gradient in Newton Methods for Linear Classification

In this section we review how CG and PCG are applied in a Newton method for linear classification.

2.1. Truncated Newton Methods

Consider training data (y_i, \mathbf{x}_i) , $i = 1, \dots, l$, where $y_i = \pm 1$ is the label and $\mathbf{x}_i \in \mathbb{R}^n$ is a feature vector. In linear classification, the model vector $\mathbf{w} \in \mathbb{R}^n$ is obtained by solving

$$\min_{\mathbf{w}} f(\mathbf{w}), \text{ where } f(\mathbf{w}) \equiv \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \xi(y_i \mathbf{w}^T \mathbf{x}_i).$$

In $f(\mathbf{w})$, $\mathbf{w}^T \mathbf{w} / 2$ is the L2-regularization term, $\xi(y_i \mathbf{w}^T \mathbf{x}_i)$ is the loss function and $C > 0$ is a parameter to balance the two terms. Here we consider logistic and L2 losses

$$\xi_{\text{LR}}(y \mathbf{w}^T \mathbf{x}) = \log(1 + \exp(-y \mathbf{w}^T \mathbf{x})), \quad (2)$$

$$\xi_{\text{L2}}(y \mathbf{w}^T \mathbf{x}) = (\max(0, 1 - y \mathbf{w}^T \mathbf{x}))^2. \quad (3)$$

Truncated Newton methods are one of the main approaches for large-scale linear classification. It iteratively finds a direction \mathbf{s}_k that minimizes the quadratic approximation of

$$f(\mathbf{w}_k + \mathbf{s}) - f(\mathbf{w}_k) \approx q_k(\mathbf{s}) \equiv \nabla f(\mathbf{w}_k)^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \nabla^2 f(\mathbf{w}_k) \mathbf{s}, \quad (4)$$

where \mathbf{w}_k is the current iterate, $\nabla f(\mathbf{w}_k)$ and $\nabla^2 f(\mathbf{w}_k)$ are the gradient and the Hessian, respectively. Because L2 Loss is not twice differentiable, we can consider the generalized Hessian matrix (Mangasarian, 2002). The direction \mathbf{s}_k from minimizing (4) can be obtained by solving the linear system

$$\nabla^2 f(\mathbf{w}_k) \mathbf{s} = -\nabla f(\mathbf{w}_k). \quad (5)$$

Exactly solving (5) for large-scale problems is expensive and furthermore $\nabla^2 f(\mathbf{w}_k)$ may be too large to be stored. Currently, conjugate gradient (CG) methods are commonly used to approximately solve (5) for obtaining a truncated Newton direction. A CG procedure involves a sequence of Hessian-vector products. Past developments (e.g., Keerthi

and DeCoste, 2005; Lin et al., 2008) have shown that the special structure of the Hessian in (6) allows us to conduct Hessian-vector products without explicitly forming the matrix. Specifically, from $f(\mathbf{w})$ we have

$$H_k = \nabla^2 f(\mathbf{w}_k) = I + CX^TDX, \quad (6)$$

where D is a diagonal matrix with $D_{ii} = \xi''(y_i \mathbf{w}_k^T \mathbf{x}_i)$, I is the identity matrix, and $X = [\mathbf{x}_1, \dots, \mathbf{x}_l]^T$ is the data matrix. The Hessian-vector product is conducted by

$$\nabla^2 f(\mathbf{w})\mathbf{s} = (I + CX^TDX)\mathbf{s} = \mathbf{s} + CX^T(DX\mathbf{s}). \quad (7)$$

The CG procedure stops after, for example, the following relative error in solving (5) is small.

$$\|\nabla^2 f(\mathbf{w}_k)\mathbf{s} + \nabla f(\mathbf{w}_k)\| \leq \varepsilon_{CG} \|\nabla f(\mathbf{w}_k)\|, \quad (8)$$

where ε_{CG} is a small positive value.¹

To ensure the convergence, after finding a direction \mathbf{s}_k , we should adjust the step size along \mathbf{s}_k (line search strategy) or decide if \mathbf{s}_k should be accepted (trust region method). In this work we focus on the trust region approach because first, it is used in the package LIBLINEAR (Fan et al., 2008), the platform considered in this study, and second, for logistic regression, in some situations the trust region approach is superior (Hsia et al., 2017).

2.2. Trust-region Newton Method

A trust region method indirectly adjusts the step size by finding a direction \mathbf{s}_k within a trust region. The direction is taken if it results in a sufficient function-value reduction. The size of the trust region is then adjusted. Here we mainly follow the description in Hsia et al. (2017).

Given a trust region with size Δ_k at the k th iteration, we compute the approximate Newton direction \mathbf{s}_k by solving the following trust-region sub-problem:

$$\min_{\mathbf{s}} q_k(\mathbf{s}) \quad \text{subject to} \quad \|\mathbf{s}\| \leq \Delta_k, \quad (9)$$

where $q_k(\mathbf{s})$ is defined in (4). Then, we check the ratio between the real and the predicted reduction of $f(\mathbf{w})$:

$$\rho_k = \frac{f(\mathbf{w}_k + \mathbf{s}_k) - f(\mathbf{w}_k)}{q_k(\mathbf{s}_k)}. \quad (10)$$

The iterate \mathbf{w} is updated only if the ratio is large enough:

$$\mathbf{w}_{k+1} = \begin{cases} \mathbf{w}_k + \mathbf{s}_k, & \text{if } \rho > \eta_0, \\ \mathbf{w}_k, & \text{if } \rho \leq \eta_0, \end{cases} \quad (11)$$

where $\eta_0 > 0$ is a pre-defined constant. The trust-region size Δ_k is adjusted by comparing the actual and the predicted function-value reduction. More details can be found in, for example, Lin et al. (2008). We summarize a trust region Newton method in Algorithm 1.

Solving the sub-problem (9) is similar to solving the linear system (5) though a constraint $\|\mathbf{s}\| \leq \Delta_k$ must be satisfied. A classic approach in Steihaug (1983) applies CG with the initial $\mathbf{s} = \mathbf{0}$ and has that $\|\mathbf{s}\|$ is monotonically increasing. Then CG stops after either (8) is satisfied or an \mathbf{s} on the boundary of the trust region is obtained. The procedure to solve the trust-region sub-problem is given in Algorithm 2.

1. In experiments, we use $\varepsilon_{CG} = 0.1$ by following the setting in the software LIBLINEAR (Lin et al., 2008).

Algorithm 1: A CG-based trust region Newton method

```

1 Given  $\mathbf{w}_0$ . for  $k = 0, 1, 2, \dots$  do
2   |   Approximately solve trust-region sub-problem (9) by the CG method to obtain a
   |   direction  $\mathbf{s}_k$ .
3   |   Compute  $\rho_k$  via (10).
4   |   Update  $\mathbf{w}_k$  to  $\mathbf{w}_{k+1}$  according to (11).
5   |   Update  $\Delta_{k+1}$  (details not discussed).
6 end

```

Algorithm 2: CG for the trust-region sub-problem (9)

```

1 Given  $\varepsilon_{\text{CG}} < 1, \Delta_k > 0$ , let  $\bar{\mathbf{s}} = \mathbf{0}, \mathbf{r} = \mathbf{d} = -\nabla f(\mathbf{w}_k)$ 
2  $\text{rTr} \leftarrow \mathbf{r}^T \mathbf{r}$ 
3 while True do
4   |   if  $\sqrt{\text{rTr}} < \varepsilon_{\text{CG}} \|\nabla f(\mathbf{w}_k)\|$  then
5   |   |   return  $\mathbf{s}_k = \bar{\mathbf{s}}$ 
6   |   end
7   |    $\mathbf{v} \leftarrow \nabla^2 f(\mathbf{w}_k) \mathbf{d}, \alpha \leftarrow \text{rTr} / (\mathbf{d}^T \mathbf{v}), \bar{\mathbf{s}} \leftarrow \bar{\mathbf{s}} + \alpha \mathbf{d}$ 
8   |   if  $\|\bar{\mathbf{s}}\| \geq \Delta_k$  then
9   |   |    $\bar{\mathbf{s}} \leftarrow \bar{\mathbf{s}} - \alpha \mathbf{d}$ 
10  |   |   compute  $\tau$  such that  $\|\bar{\mathbf{s}} + \tau \mathbf{d}\| = \Delta_k$ 
11  |   |   return  $\mathbf{s}_k = \bar{\mathbf{s}} + \tau \mathbf{d}$ 
12  |   end
13  |    $\mathbf{r} \leftarrow \mathbf{r} - \alpha \mathbf{v}, \text{rTr}_{\text{new}} \leftarrow \mathbf{r}^T \mathbf{r}$ 
14  |    $\beta \leftarrow \text{rTr}_{\text{new}} / \text{rTr}, \mathbf{d} \leftarrow \mathbf{r} + \beta \mathbf{d}$ 
15  |    $\text{rTr} \leftarrow \text{rTr}_{\text{new}}$ 
16 end

```

We explain that the computational cost per iteration of truncated Newton methods is roughly

$$\mathcal{O}(nl) \times (\# \text{ CG steps}). \quad (12)$$

From (7), $\mathcal{O}(nl)$ is the cost per CG step. Besides CG, function and gradient evaluation costs about the same as one CG step. Because in general each Newton iteration requires several CG steps, we omit function/gradient evaluation in (12) and hence the total cost is proportional to the total number of CG steps. If X is a sparse matrix with $\#\text{nnz}$ non-zero entries, then the $\mathcal{O}(nl)$ term can be replaced by $\mathcal{O}(\#\text{nnz})$.

2.3. PCG for Trust-region Sub-problem

From (12), reducing the number of CG steps can speed up the Newton method. The preconditioning technique (Concus et al., 1976) has been widely used to possibly improve the condition of linear systems and reduce the number of CG steps. PCG considers a preconditioner

$$M = EE^T \approx \nabla^2 f(\mathbf{w}_k)$$

Algorithm 3: PCG for solving the transformed trust-region sub-problem (14). Assume M has not been factorized to EE^T .

```

1 Given  $\varepsilon_{\text{CG}} < 1, \Delta_k > 0$ 
2 Let  $\bar{\mathbf{s}} = \mathbf{0}, \mathbf{r} = -\nabla f(\mathbf{w}_k), \mathbf{d} = \mathbf{z} = M^{-1}\mathbf{r}, \gamma = \mathbf{r}^T \mathbf{z}$ 
3 while True do
4   if  $\sqrt{\mathbf{r}^T \mathbf{z}} < \varepsilon_{\text{CG}} \|\nabla f(\mathbf{w}_k)\|$  then
5     return  $\mathbf{s}_k = \bar{\mathbf{s}}$ 
6   end
7    $\mathbf{v} \leftarrow \nabla^2 f(\mathbf{w}_k) \mathbf{d}, \alpha \leftarrow \mathbf{r}^T \mathbf{z} / (\mathbf{d}^T \mathbf{v})$ 
8    $\bar{\mathbf{s}} \leftarrow \bar{\mathbf{s}} + \alpha \mathbf{d}$ 
9   if  $\|\bar{\mathbf{s}}\|_M \geq \Delta_k$  then
10     $\bar{\mathbf{s}} \leftarrow \bar{\mathbf{s}} - \alpha \mathbf{d}$ 
11    compute  $\tau$  such that  $\|\bar{\mathbf{s}} + \tau \mathbf{d}\|_M = \Delta_k$ 
12    return  $\mathbf{s}_k = \bar{\mathbf{s}} + \tau \mathbf{d}$ 
13  end
14   $\mathbf{r} \leftarrow \mathbf{r} - \alpha \mathbf{v}, \mathbf{z} \leftarrow M^{-1} \mathbf{r}$ 
15   $\gamma^{\text{new}} \leftarrow \mathbf{r}^T \mathbf{z}, \beta \leftarrow \gamma^{\text{new}} / \gamma$ 
16   $\mathbf{d} \leftarrow \mathbf{z} + \beta \mathbf{d}$ 
17   $\gamma \leftarrow \gamma^{\text{new}}$ 
18 end

```

and transforms the original linear system (5) to

$$E^{-1} \nabla^2 f(\mathbf{w}_k) E^{-T} \hat{\mathbf{s}} = -E^{-1} \nabla f(\mathbf{w}_k). \quad (13)$$

If the approximation is good, the condition number of $E^{-1} \nabla^2 f(\mathbf{w}_k) E^{-T}$ is close to 1 and less CG steps are needed. After obtaining $\hat{\mathbf{s}}_k$, we can get the original solution by $\mathbf{s}_k = E^{-T} \hat{\mathbf{s}}_k$. More discussion about the condition number and the convergence of CG is in the supplement.

To apply PCG for the trust-region sub-problem, we follow [Steihaug \(1983\)](#) to modify the sub-problem (9) to

$$\begin{aligned} \min_{\hat{\mathbf{s}}} \quad & \frac{1}{2} \hat{\mathbf{s}}^T (E^{-1} \nabla^2 f(\mathbf{w}_k) E^{-T}) \hat{\mathbf{s}} + (E^{-1} \nabla f(\mathbf{w}_k))^T \hat{\mathbf{s}} \\ \text{subject to} \quad & \|\hat{\mathbf{s}}\| \leq \Delta_k. \end{aligned} \quad (14)$$

Then the same procedure in Algorithm 2 can be applied to solve the new sub-problem. See details in Algorithm I in supplementary materials. By comparing Algorithms 2 and I, clearly the extra cost of PCG is for calculating the product between E^{-1} (or E^{-T}) and a vector; see

$$E^{-1} (\nabla^2 f(\mathbf{w}_k) (E^{-T} \hat{\mathbf{d}})) \quad (15)$$

in Algorithm I. In some situations the factorization of $M = EE^T$ is not practically viable, but it is known ([Golub and Van Loan, 1996](#)) that PCG can be performed without a factorization of M . We show the procedure in Algorithm 3. The extra cost of PCG is shifted from (15) to the product between M^{-1} and a vector; see

$$M^{-1} \mathbf{r}$$

in Algorithm 3.

If the factorization $M = EE^T$ is available, a practical issue is whether Algorithm 3 or Algorithm I in supplementary materials should be used. We give a detailed investigation in Section IV of supplementary materials. In the end we choose Algorithm 3 for our implementation.

3. Existing Preconditioners

We discuss several existing preconditioners which have been applied to linear classification.

3.1. Diagonal Preconditioner

It is well known that extracting all diagonal elements in the Hessian can form a simple preconditioner.

$$M = \text{diag}(H_k), \text{ where } M_{ij} = \begin{cases} (H_k)_{ij}, & \text{if } i = j, \\ 0, & \text{otherwise.} \end{cases}$$

From (6),

$$(H_k)_{jj} = 1 + C \sum_i D_{ii} X_{ij}^2,$$

so constructing the preconditioner costs $\mathcal{O}(nl)$. At each CG step, the cost of computing $M^{-1}\mathbf{r}$ is $\mathcal{O}(n)$. Thus, the extra cost of using the diagonal preconditioner is

$$\mathcal{O}(n) \times (\# \text{ CG steps}) + \mathcal{O}(nl). \tag{16}$$

In comparison with (12), the cost of using a diagonal preconditioner at each Newton iteration is insignificant.

For linear classification, Lin et al. (2008) have applied the diagonal preconditioner, but conclude that it is not consistently better. However, Chin et al. (2016) pointed out that the comparison is conducted by strictly solving the optimization problem. Because in general a less accurate optimization solution is enough for machine learning applications, Chin et al. (2016) find that diagonal preconditioning is practically useful. We will conduct a thorough evaluation in Section 5.1.

3.2. Subsampled Hessian as Preconditioner

In Section 2.3, we have shown that a good preconditioner M should satisfy $M \approx \nabla^2 f(\mathbf{w}_k)$. In linear classification, the object function involves the sum of training losses. If we randomly select \bar{l} instance-label pairs $(\bar{y}_i, \bar{\mathbf{x}}_i)$ and construct a subsampled Hessian (Byrd et al., 2011)

$$\bar{H}_k = I + C \frac{\bar{l}}{\bar{l}} \bar{X}^T \bar{D} \bar{X},$$

where $\bar{X} = [\bar{\mathbf{x}}_1^T, \dots, \bar{\mathbf{x}}_{\bar{l}}^T] \in \mathbb{R}^{\bar{l} \times n}$ and $\bar{D} \in \mathbb{R}^{\bar{l} \times \bar{l}}$ is a diagonal matrix with $\bar{D}_{ii} = \xi''(\bar{y}_i \mathbf{w}_k^T \bar{\mathbf{x}}_i)$, then \bar{H}_k is an unbiased estimator for $\nabla^2 f(\mathbf{w}_k)$. Therefore, $M = \bar{H}_k$ can be considered as a reasonable preconditioner.

However, $\bar{H}_k \in \mathbb{R}^{n \times n}$ has the same size as H_k so it may be too large to be stored. To get $\bar{H}_k^{-1} \mathbf{r}$ in PCG, [Ma and Takáč \(2016\)](#) apply the following Woodbury formula. Consider $A \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{n \times m}$ and $V \in \mathbb{R}^{m \times n}$. If A and $(I + VA^{-1}U)$ are invertible, then

$$(A + UV)^{-1} = A^{-1} - A^{-1}U(I + VA^{-1}U)^{-1}VA^{-1}.$$

To apply this formula on the sub-sampled Hessian \bar{H}_k , we first let

$$G = \left(\frac{lC}{\bar{l}}D\right)^{\frac{1}{2}}.$$

Then

$$\bar{H}_k = I_n + (\bar{X}^T G)I_{\bar{l}}(G\bar{X})$$

and we have

$$\bar{H}_k^{-1} = I_n - \bar{X}^T G(I_{\bar{l}} + G\bar{X}\bar{X}^T G)^{-1}G\bar{X}, \quad (17)$$

where $G = \left(\frac{lC}{\bar{l}}D\right)^{\frac{1}{2}}$. If the chosen \bar{l} satisfies $\bar{l} \ll l$, then

$$(I_{\bar{l}} + G\bar{X}\bar{X}^T G)^{-1} \in \mathbb{R}^{\bar{l} \times \bar{l}}$$

can be calculated in $\mathcal{O}(n\bar{l}^2 + \bar{l}^3)$ time and is small enough to be stored. Note that \bar{l}^3 is for calculating the inverse. Then in PCG the product

$$\bar{H}_k^{-1} \mathbf{r} = \mathbf{r} - \bar{X}^T (G((I_{\bar{l}} + G\bar{X}\bar{X}^T G)^{-1}(G(\bar{X}\mathbf{r}))))$$

can be done in the cost of $\mathcal{O}(n\bar{l}) + \mathcal{O}(\bar{l}^2)$. For large and sparse data, $\bar{l} \ll n$, so at each Newton iteration, the additional cost of using \bar{H}_k as the preconditioner is

$$\mathcal{O}(n\bar{l}) \times (\# \text{ CG steps}) + \mathcal{O}(n\bar{l}^2). \quad (18)$$

A comparison with (12) shows that we can afford this cost if $\bar{l} \ll l$. Further, from (16) and (18), using a sub-sampled Hessian as the preconditioner is in general more expensive than diagonal preconditioning. We will discuss the selection of \bar{l} in Section 5.2.

4. Our Proposed Method

It is hard to find a preconditioner that works well in all cases. We give a simple example showing that the diagonal preconditioner may not reduce the condition number. If

$$H_k = I + \begin{bmatrix} 2 & 2 & 2 \\ 2 & 3 & 3 \\ 2 & 3 & 4 \end{bmatrix} \text{ and } M = EE^T = \text{diag}(H_k),$$

then

$$\kappa(E^{-1}H_kE^{-T}) = 6.781 > \kappa(H_k) = 6.723.$$

In a Newton method, each iteration involves a sub-problem. A preconditioner may be useful for some sub-problems, but not for others. Thus it is difficult to ensure the shorter overall training time of using PCG. We find that if some Newton iterations need many CG

steps because of inappropriate preconditioning, the savings by PCG in subsequent Newton iterations may not be able to recover the loss. Thus we check if a more robust setting can be adopted. Assume

$$M = EE^T \approx H_k$$

is any preconditioner considered for the current sub-problem. We hope to derive a new preconditioner $\bar{M} = \bar{E}\bar{E}^T$ that satisfies

$$\kappa(\bar{E}^{-1}H_k\bar{E}^{-T}) \approx \min\{\kappa(H_k), \kappa(E^{-1}H_kE^{-T})\}. \quad (19)$$

This property avoids the situation when $E^{-1}H_kE^{-T}$ has a larger condition number and somehow ensures that PCG is useful at each Newton iteration. Note that we assume the existence of the factorization EE^T and $\bar{E}\bar{E}^T$ only for discussing theoretical properties (19). In practice we can just use M or \bar{M} as indicated in Section 2.3.

The next issue is how to achieve the inequality (19). We offer two approaches.

4.1. Running CG and PCG in Parallel

If in an ideal situation we can run CG and PCG in parallel, then we can conjecture that the one needing fewer steps has a smaller condition number. Therefore, the following rule can choose between CG and PCG at each Newton iteration.

$$\bar{M} = \begin{cases} I, & \text{if CG uses less steps,} \\ M, & \text{if PCG uses less steps.} \end{cases}$$

Because we consider preconditioners that do not incur much extra cost, this setting is similar to checking which one finishes first. Our proposed strategy can be easily implemented in a multi-core or distributed environment.

4.2. Weighted Average of the Preconditioner and the Identity Matrix

The method in Section 4.1 to run CG and PCG in parallel is a direct way to achieve (19), but it is not useful in a single-thread situation. Unfortunately, developing a new preconditioner satisfying (19) is extremely difficult. In general we do not have that one preconditioner is consistently better than another. Therefore, we consider a more modest goal of improving the situation when preconditioning is harmful. The property (19) is revised to

$$\kappa(\bar{E}^{-1}H_k\bar{E}^{-T}) \not\approx \max\{\kappa(H_k), \kappa(E^{-1}H_kE^{-T})\},$$

or equivalently

$$\kappa(\bar{E}^{-1}H_k\bar{E}^{-T}) < \max\{\kappa(H_k), \kappa(E^{-1}H_kE^{-T})\}. \quad (20)$$

To achieve (20), we prove in the following theorem that one possible setting is by a weighted average of M and an identity matrix (details of the proof are in supplementary materials).

Theorem 1 *For any symmetric positive definite H_k and any M with $\kappa(M^{-1}H_k) \neq \kappa(H_k)$, the following matrix satisfies (20) for all $0 < \alpha < 1$.*

$$\bar{M} = \alpha M + (1 - \alpha)I.$$

Data sets	#instances	#features	$\log_2(C_{\text{Best}})$	
			LR	L2
news20	19,996	1,355,191	9	3
yahookr	460,554	3,052,939	6	1
url	2,396,130	3,231,962	-7	-10
kddb	19,264,097	29,890,095	-1	-4
criteo	45,840,617	1,000,000	-15	-12
kdd12	149,639,105	54,686,452	-4	-11

Table 1: Data statistics. C_{Best} is the regularization parameter selected by cross validation.

Note that $\alpha = 0$ and 1 respectively indicate that CG (no preconditioning) and PCG (with preconditioner M) are used. By choosing an $\alpha \in (0, 1)$, we avoid the worse one of the two settings.

If M is the diagonal preconditioner in Section 3.1, then

$$\bar{M} = \alpha \times \text{diag}(H_k) + (1 - \alpha) \times I. \quad (21)$$

In Section 5.1 we will investigate the effectiveness of using (21) and the selection of α .

5. Experiments

In Section 5.1, we compare the diagonal preconditioner and the proposed methods. For the subsampled-Hessian preconditioner, we check the performance in Section 5.2. Finally, an overall comparison is in Section 5.3.

We consider binary classification data sets listed in Table 1. All data sets except `yahookr` can be downloaded from [LIBSVM Data Sets \(2007\)](#). We extend the trust region Newton implementation in a popular linear classification package LIBLINEAR (version 2.11) to incorporate various preconditioners. By extending from CG to PCG, many implementation details must be addressed. See Section IV of supplementary materials.

We present results of using the logistic loss, while leaving results of the l2 loss in the supplementary materials. Some preliminary results on linear support vector regression are also given there.

In experiments, if the same type of preconditioners are compared, we check the total number of CG steps when the optimization procedure reaches the following default stopping condition of LIBLINEAR

$$\|\nabla f(\mathbf{w}_k)\| \leq \epsilon \frac{\min(\#\text{pos}, \#\text{neg})}{l} \|\nabla f(\mathbf{w}_0)\|, \quad (22)$$

where $\epsilon = 10^{-2}$, $\#\text{pos}$, $\#\text{neg}$ are the numbers of positive- and negative-labeled instances respectively, and l is the total number of instances.

From (12), if the cost of preconditioning is insignificant, a comparison on CG steps is the same as a running time comparison. However, for experiments in Section 5.3 using expensive preconditioners, we give timing results. For the regularization parameter C , we consider

$$C = C_{\text{Best}} \times \{0.01, 0.1, 1, 10, 100\}, \quad (23)$$

where C_{Best} for each data set is the value leading to the best cross validation accuracy.

(a) $C = C_{\text{Best}}$				(b) $C = 100C_{\text{Best}}$			
Data	Diag	CG or Diag	Mixed	Data	Diag	CG or Diag	Mixed
news20	1.61	1.06	0.98	news20	2.38	0.85	0.98
url	1.25	0.86	0.87	url	1.14	0.50	1.29
yahookr	0.29	0.44	0.67	yahookr	0.31	0.11	0.16
kddb	0.24	0.25	0.28	kddb	0.04	0.03	0.05
kdd12	0.19	0.19	0.31	kdd12	0.29	0.10	0.38
criteo	0.65	0.68	0.70	criteo	0.82	0.37	0.49

Table 2: The ratio between the total number of CG steps of a method and that of using standard CG. The smaller the ratio is better. Ratios larger than one are boldfaced, indicating that preconditioning is not helpful. We run the Newton method until the stopping condition (22) is satisfied. Logistic loss is used.

5.1. Diagonal Preconditioner and the Proposed Method in Section 4

We compare the following settings.

- **CG**: LIBLINEAR version 2.11, a trust-region method without preconditioning.
- **Diag**: the diagonal preconditioner in Section 3.1.
- **CG or Diag**: the technique in Section 4.1 to run CG and diagonal PCG in parallel.²
- **Mixed**: the preconditioner in (21) with $\alpha = 10^{-2}$.

In Table 2, we present

$$\frac{\text{Total \#CG steps of a method}}{\text{Total \#CG steps without preconditioner}} \quad (24)$$

by using $C = C_{\text{Best}} \times \{1, 100\}$. Thus a value smaller than 1 indicates that PCG reduces the number of CG steps. Full results under all C values listed in (23) are in supplementary materials.

From Table 2, in general the ratio of using the diagonal preconditioner is less than one, indicating that applying this preconditioner can reduce the total number of CG steps. However, for `news20` and `url`, diagonal preconditioning is not useful. We conduct a detailed investigation by checking the relationship between the accumulated number of CG steps and the following relative function-value reduction

$$\frac{f(\mathbf{w}_k) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)},$$

where \mathbf{w}^* is an approximate optimal solution by running many iterations. Figure 1 shows that the final convergence of using diagonal preconditioning is much slower than without preconditioning. Thus the diagonal preconditioner is not robust for practical use.

For the proposed approaches **CG or Diag** and **Mixed**, they are generally better or as good as the two original settings: **CG** (no preconditioning) and **Diag** (diagonal preconditioning). Therefore, the proposed techniques effectively improve the robustness of the diagonal preconditioner. In supplementary materials we provide more results including the sensitivity of the α value in the **Mixed** approach.

2. Because of checking the total number of CG steps, we can easily use a single machine to simulate the parallel setting.

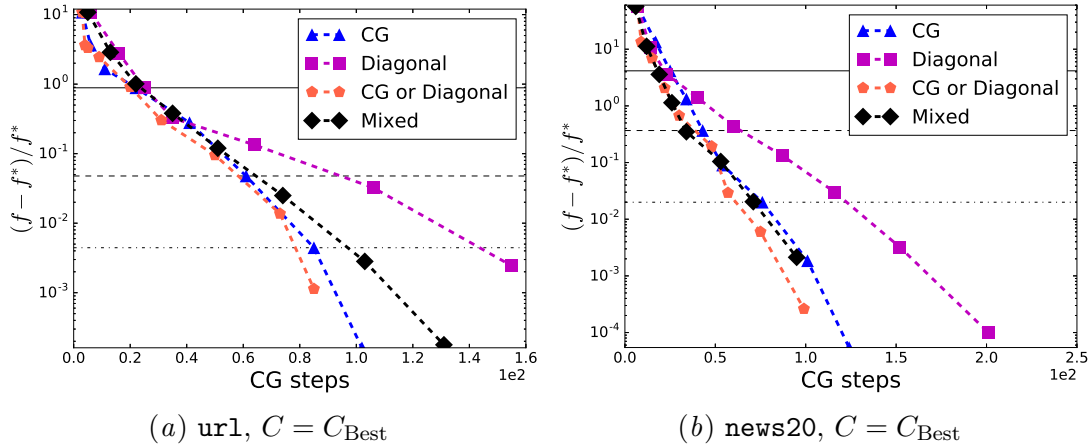


Figure 1: Comparison of the convergence of CG, Diag, CG or Diag and Mixed. Logistic loss is used. We show a relative difference to the optimal function value (log-scaled) versus the total CG steps. Horizontal lines show that LIBLINEAR’s stopping condition (22) with tolerances 10^{-1} , 10^{-2} (default), and 10^{-3} is reached; such information indicates when the training algorithm should stop.

Data	(a) $C = C_{\text{Best}}$			Data	(b) $C = 100C_{\text{Best}}$		
	SH-100	SH-1000	SH-3000		SH-100	SH-1000	SH-3000
news20	1.00	1.20	2.02	news20	0.98	1.20	1.93
url	0.92	0.73	0.72	url	1.63	1.12	1.05
yahookr	1.17	0.71	0.50	yahookr	1.01	0.71	1.11
kddb	0.98	0.50	0.35	kddb	0.89	0.69	0.38
kdd12	0.79	0.69	0.61	kdd12	1.31	1.50	1.46
criteo	0.77	0.47	0.28	criteo	1.14	0.53	0.53

Table 3: A comparison of using different \bar{l} values in the sub-sampled Hessian preconditioner. We present the ratio between the total number of CG steps of a method and that of using standard CG. The smaller the ratio is better. Ratios larger than one are boldfaced, indicating that preconditioning is not helpful. Other settings are the same as Table 2.

5.2. Subsampled Hessian as Preconditioner

We investigate the size \bar{l} in the subsampled Hessian preconditioner described in Section 3.2 by considering the following settings.

- SH-100: method in Section 3.2 with $\bar{l} = 100$.
- SH-1000: method in Section 3.2 with $\bar{l} = 1,000$.
- SH-3000: method in Section 3.2 with $\bar{l} = 3,000$.

For each setting we compare the number of CG steps with that of standard CG by calculating the ratio in (24). Note that from (18), the cost of using subsampled Hessian is high, so we will present timing results in Section 5.3.

From Table 3, a larger \bar{l} generally leads to a smaller number of CG steps. However, in some situations (e.g., news20), a larger \bar{l} is not useful. From a detailed investigation

(a) $C = C_{\text{Best}}$				(b) $C = 100C_{\text{Best}}$			
Data	Diagonal	Mixed	SH-3000	Data	Diagonal	Mixed	SH-3000
news20	1.76	1.13	44.36	news20	2.51	1.15	48.20
url	1.24	0.91	1.16	url	1.18	1.28	1.28
yahookr	0.35	0.73	1.17	yahookr	0.33	0.19	2.22
kddb	0.28	0.31	0.41	kddb	0.05	0.05	0.42
kdd12	0.15	0.22	0.37	kdd12	0.29	0.36	1.27
criteo	0.74	0.80	0.43	criteo	0.75	0.47	0.55

Table 4: Running time comparison of using different preconditioners. We show the ratio between the running time of a method and that of using standard CG. The smaller the ratio is better. Ratios larger than one are boldfaced, indicating that preconditioning is not helpful. Other settings are the same as Table 2.

in supplementary materials, we even find that the performance is sensitive to the random seeds in constructing the sub-sampled Hessian. Thus in some cases the current \bar{l} is not large enough to give a good approximation of $\nabla^2 f(\mathbf{w}_k)$. Our observation is slightly different from that in [Ma and Takáč \(2016\)](#), which uses only $\bar{l} = 100$. It is unclear why different results are reported, but this situation seems to indicate that selecting a suitable \bar{l} is not an easy task.

5.3. Running-time Comparison of Different Preconditioners

We conduct an overall timing comparison. Results in Table 4 lead to the following observations.

1. In some situations, the cost of using a preconditioner is not negligible. For `news20`, from Table 3, SH-3000 is not much worse in terms of CG steps, but is dramatically slower in terms of time.
2. Overall `Mixed` is the best approach. It is more robust than `Diag`, and is often much faster than `CG` and `SH-3000`.

For a further illustration, we present a detailed comparison in Figure 2, where the convergences in terms of both the number of CG steps and the running time are checked. Clearly, `CG` and `Diag` are not very robust. They have the fastest and the slowest convergences on different problems. We also see that a good reduction on the number of CG steps by `SH-3000` may not lead to shorter training time. A complete set of figures by using all data sets is in supplementary materials.

6. Conclusions

In this work we show that applying preconditioners in Newton methods for linear classification is not an easy task. Improvements made at one Newton iteration may not lead to better overall convergence. We propose using a reliable preconditioner at each iteration. The idea is that between the setting of no preconditioning and the setting of using one particular preconditioner, we try to select the better one. If the selection is not possible in practice, we propose using a weighted combination to ensure that at least the worse one is not considered. Experiments confirm that the proposed method leads to faster overall

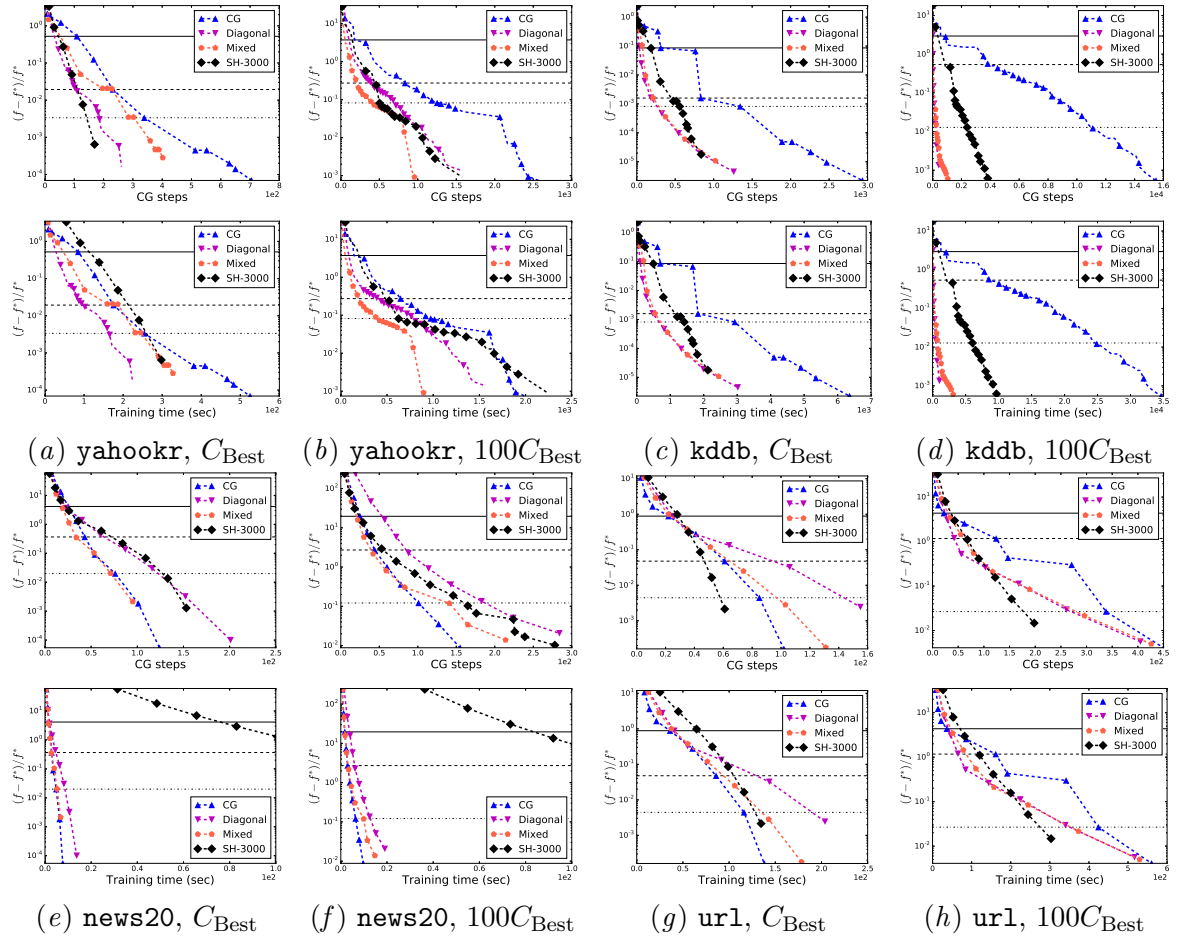


Figure 2: Convergence of using different preconditioners. We present the results of using $C = C_{\text{Best}}$ and $C = 100C_{\text{Best}}$. For the same data set and under the same C value, the x -axis of the upper figure is the cumulative number of CG steps and the x -axis of the lower one is the running time. We align curves of the approach CG in upper and lower figures for an easy comparison. Other settings are the same as those in Figure 1.

convergence. A corresponding implementation has been included in a package for public use.

References

Richard H. Byrd, Gillian M. Chin, Will Neveitt, and Jorge Nocedal. On the use of stochastic Hessian information in optimization methods for machine learning. *SIAM Journal on Optimization*, 21(3):977–995, 2011.

Wei-Sheng Chin, Bo-Wen Yuan, Meng-Yuan Yang, and Chih-Jen Lin. A note on diagonal preconditioner in large-scale logistic regression. Technical report, Department of

- Computer Science and Information Engineering, National Taiwan University, 2016. URL <https://www.csie.ntu.edu.tw/~cjlin/papers/logistic/pcgnote.pdf>.
- Paul Concus, Gene H Golub, and Dianne P O’Leary. A generalized conjugate gradient method for the numerical solution of elliptic partial differential equations. Technical report, 1976.
- 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>.
- G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, third edition, 1996.
- Chih-Yang Hsia, Ya Zhu, and Chih-Jen Lin. A study on trust region update rules in Newton methods for large-scale linear classification. In *Proceedings of the Asian Conference on Machine Learning (ACML)*, 2017. URL <http://www.csie.ntu.edu.tw/~cjlin/papers/newtron/newtron.pdf>.
- 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.
- LIBSVM Data Sets, 2007. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>.
- 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>.
- Chenxin Ma and Martin Takáč. Distributed inexact damped Newton method: data partitioning and load-balancing. *arXiv preprint arXiv:1603.05191*, 2016.
- Olvi L. Mangasarian. A finite Newton method for classification. *Optimization Methods and Software*, 17(5):913–929, 2002.
- Stephen G Nash. Preconditioning of truncated-Newton methods. *SIAM Journal on Scientific and Statistical Computing*, 6(3):599–616, 1985.
- Trond Steihaug. The conjugate gradient method and trust regions in large scale optimization. *SIAM Journal on Numerical Analysis*, 20:626–637, 1983.
- Yuchen Zhang and Xiao Lin. DiSCO: Distributed optimization for self-concordant empirical loss. In *Proceedings of the Thirty Second International Conference on Machine Learning (ICML)*, pages 362–370, 2015.