

Supplementary Materials for “Subsampled Hessian Newton Methods for Supervised Learning”

Chien-Chih Wang, Chun-Heng Huang, Chih-Jen Lin

Department of Computer Science, National Taiwan University, Taipei 10617, Taiwan

I Introduction

This document presents some materials that are not included in the paper. In Section II, we show experiments of using a subset of data for calculating the gradient. In Section III, we discuss some variants of the proposed method. Section IV gives details of applying our method to maximum entropy.

II Using a Subset of Data for Calculating the Gradient

We mentioned in Section 1 that Byrd et al. [2011] consider using a subset R so that

$$\nabla f(\mathbf{w}) \approx \frac{1}{|R|} \sum_{i \in R} \nabla \xi(\mathbf{w}; \mathbf{x}_i, y_i).$$

Byrd et al. [2011] conducted experiments by using $R = \{1, \dots, l\}$, so only the Hessian is approximated by a subset of data. We follow the same setting in the paper. One reason of not subsampling points for gradient evaluation is that in the line search procedure we still need to access the whole set for computing the function value.

Here we compare the following two settings:

1. *Method 2*: the proposed method in the paper; see Section 6.1.
2. *Method 2-sg*: the method is the same as *Method 2* except that data are subsampled for calculating the gradient. For example, *Method 2-sg-1/2-CG10* is that we only use a subset R with $|R| = 50\%l$ to derive the gradient.

Note that a subset S_k of R_k is further selected for the Hessian Calculation. The comparison results on logistic regression are in Figure II.1. Clearly, *Method 2-sg* is much slower. Our results indicate that because one pass of data can yield both function and gradient values, there may be no need to subsample points for the gradient calculation.

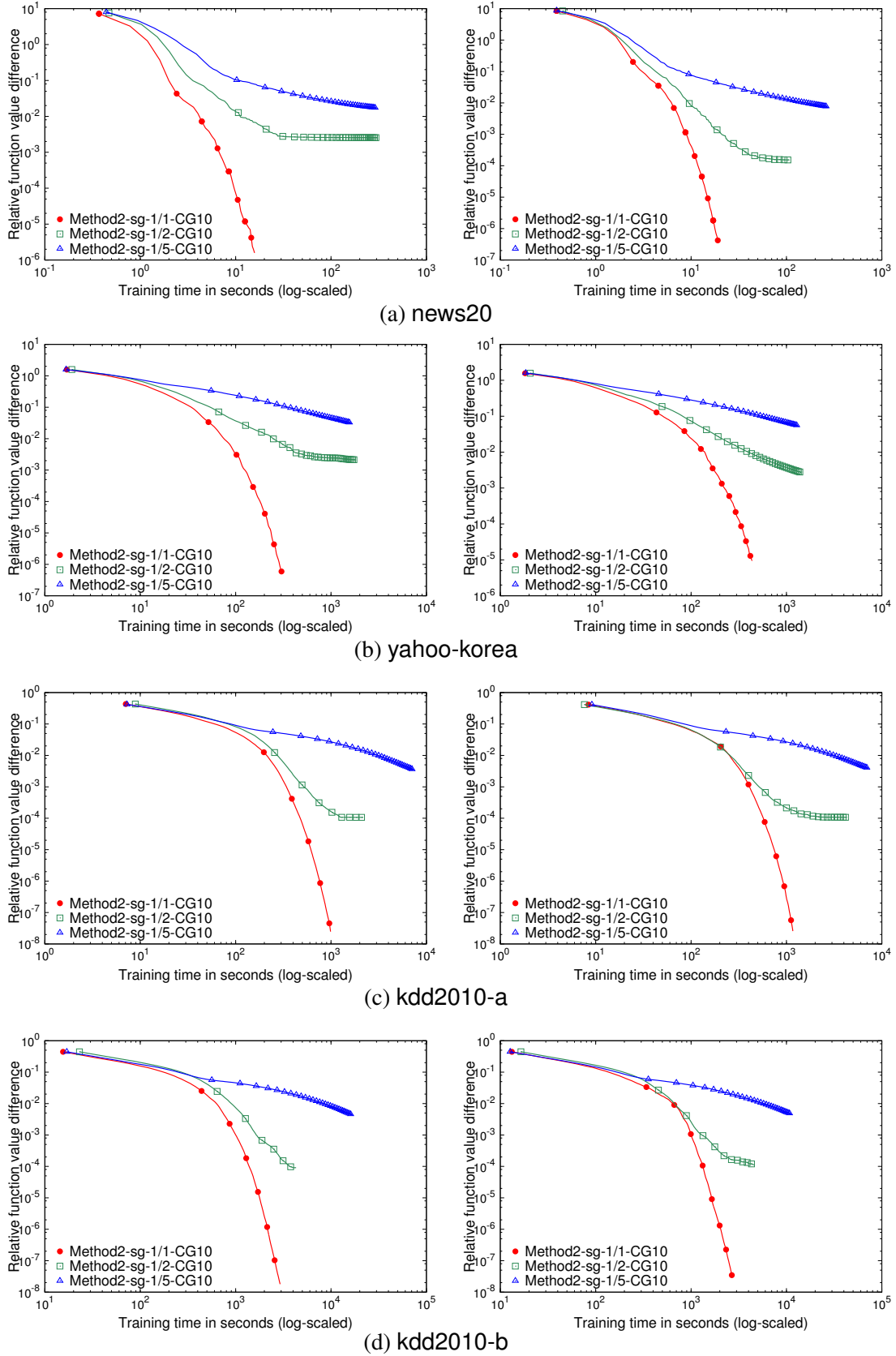


Figure II.1: Experiments on approaches with/without using a subset for gradient calculation. Logistic regression is considered. We present running time (in seconds) versus the relative difference to the optimal function value. Both x -axis and y -axis are log-scaled. Left: $|S_k|/l = 5\%$. Right: $|S_k|/l = 1\%$.

III Some Variants of the Proposed Method

In Section III.1, we discuss the selection of $\bar{\mathbf{d}}_k$. In Section III.2, we investigate the use of non-negative β_1 and β_2 for calculating the direction $\beta_1 \mathbf{d}_k + \beta_2 \bar{\mathbf{d}}_k$ in our proposed method. Experimental comparisons are in Section III.3.

III.1 Selection of $\bar{\mathbf{d}}_k$

An important issue for *Method 2* is how to select an appropriate $\bar{\mathbf{d}}_k$. The choice of $\bar{\mathbf{d}}_k$ affects the convergence speed. In the paper, we use $\bar{\mathbf{d}}_k = \mathbf{d}_{k-1}$, $k \geq 1$. Here we try another setting

$$\bar{\mathbf{d}}_k = -\nabla f(\mathbf{w}_k).$$

III.2 Using Non-Negative β_1 and β_2

The coefficients β_1 and β_2 of solving (13) may be negative. It is interesting to check if imposing non-negativity can lead to a better direction. Therefore, we replace (13) with the following optimization problem.

$$\begin{aligned} \min_{\beta_1, \beta_2} \quad & \frac{1}{2} (\beta_1 \mathbf{d}_k + \beta_2 \bar{\mathbf{d}}_k)^T H_k (\beta_1 \mathbf{d}_k + \beta_2 \bar{\mathbf{d}}_k) + \nabla f(\mathbf{w}_k)^T (\beta_1 \mathbf{d}_k + \beta_2 \bar{\mathbf{d}}_k) \\ \text{subject to} \quad & \beta_1 \geq 0, \beta_2 \geq 0. \end{aligned} \quad (\text{III.1})$$

Although (13) has a closed-form solution, (III.1) does not. We derive a solution procedure by checking its optimality condition. Let

$$\begin{aligned} a &= \mathbf{d}_k^T H_k \mathbf{d}_k, \quad b = \bar{\mathbf{d}}_k^T H_k \mathbf{d}_k, \quad c = \bar{\mathbf{d}}_k^T H_k \bar{\mathbf{d}}_k, \\ e &= -\nabla f(\mathbf{w}_k)^T \mathbf{d}_k, \quad f = -\nabla f(\mathbf{w}_k)^T \bar{\mathbf{d}}_k. \end{aligned}$$

Problem (III.1) can be rewritten as

$$\begin{aligned} \min_{\beta_1, \beta_2} \quad & \frac{1}{2} [\beta_1 \quad \beta_2] \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} - [e \quad f] \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \\ \text{subject to} \quad & \beta_1 \geq 0, \beta_2 \geq 0. \end{aligned} \quad (\text{III.2})$$

We show that if

$$\mathbf{d}_k \neq \mathbf{0}, \quad \bar{\mathbf{d}}_k \neq \mathbf{0}, \quad \text{and} \quad \mathbf{d}_k \neq \bar{\mathbf{d}}_k, \quad (\text{III.3})$$

then

$$a > 0, \quad c > 0, \quad \text{and} \quad ac - b^2 > 0. \quad (\text{III.4})$$

The first two inequalities hold because H_k is positive definite. For the third inequality, we have Cholesky factorization of H_k

$$H_k = LL^T,$$

where L is lower triangular and invertible. Then

$$ac - b^2 = \|L\mathbf{d}_k\|^2 \|L\bar{\mathbf{d}}_k\|^2 - ((L\mathbf{d}_k)^T (L\bar{\mathbf{d}}_k))^2 > 0$$

by Cauchy inequality and the assumption $\mathbf{d}_k \neq \bar{\mathbf{d}}_k$. We check if the three conditions in (III.3) can be easily fulfilled. Because \mathbf{w}_k is not an optimal solution yet, $\nabla f(\mathbf{w}_k) \neq \mathbf{0}$ and so is \mathbf{d}_k from the CG procedure of using $\nabla f(\mathbf{w}_k)$ at the first step. The direction $\bar{\mathbf{d}}_k$ can be easily chosen so that the other two conditions hold.

The KKT optimality condition of (III.2) is that there are λ_1 and λ_2 such that

$$\begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} - \begin{bmatrix} e \\ f \end{bmatrix} = \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix},$$

$$\lambda_1 \beta_1 = 0, \lambda_2 \beta_2 = 0,$$

$$\beta_1 \geq 0, \beta_2 \geq 0, \lambda_1 \geq 0, \lambda_2 \geq 0.$$

We consider the following three situations

$$ec - bf \geq 0, af - be \geq 0 \quad (\text{III.5})$$

$$bf - ce \geq 0, f \geq 0 \quad (\text{III.6})$$

$$be - af \geq 0, e \geq 0 \quad (\text{III.7})$$

- If (III.5) holds,

$$\beta_1 = \frac{ec - bf}{ac - b^2}, \beta_2 = \frac{af - be}{ac - b^2}, \lambda_1 = 0, \lambda_2 = 0$$

satisfy the KKT condition.

- If (III.6) holds,

$$\beta_1 = 0, \beta_2 = \frac{f}{c}, \lambda_1 = \frac{bf}{c} - e, \lambda_2 = 0$$

satisfy the KKT condition.

- If (III.7) holds,

$$\beta_1 = \frac{e}{a}, \beta_2 = 0, \lambda_1 = 0, \lambda_2 = \frac{be}{a} - f$$

satisfy the KKT condition.

The following theorem shows that for all remaining situations, $(\beta_1, \beta_2) = (0, 0)$ is an optimal solution.

Theorem 1. *If none of (III.5), (III.6), (III.7) holds, then*

$$\beta_1 = \beta_2 = 0$$

is optimal for (III.1).

Proof. We show that if none of (III.5), (III.6), (III.7) holds, then

$$f \leq 0 \text{ and } e \leq 0. \quad (\text{III.8})$$

Then

$$\beta_1 = 0, \beta_2 = 0, \lambda_1 = -e, \lambda_2 = -f$$

satisfy the KKT condition.

If none of (III.5), (III.6), (III.7) holds, then we have

$$\begin{aligned} & (ec - bf < 0 \text{ or } af - be < 0), \text{ and} \\ & (bf - ce < 0 \text{ or } f < 0), \text{ and} \\ & (be - af < 0 \text{ or } e < 0). \end{aligned} \tag{III.9}$$

We argue that the above condition implies (III.8). Otherwise, if (III.8) does not hold, we have

$$f > 0 \text{ or } e > 0. \tag{III.10}$$

Consider the first situation where

$$f > 0.$$

Then (III.9) implies

$$\begin{aligned} bf - ce &< 0, \\ af - be &< 0, \\ e &< 0. \end{aligned}$$

From (III.4),

$$b < \frac{ce}{f} < 0 \text{ and } b < \frac{af}{e} < 0$$

lead to

$$b^2 > ac,$$

which is a contradiction to (III.4). The situation for

$$e > 0$$

is similar. Therefore, (III.8) holds and the proof is complete. \square

III.3 Experiments

We compare the following three methods in Figures III.2 and III.3 respectively for logistic regression and l2-loss SVM.

1. *Method 2*: the method proposed in the paper.
2. *Method 2-g*: the same as *Method 2* except using $\bar{\mathbf{d}}_k = -\nabla f(\mathbf{w}_k)$.
3. *Method 2-con*: the same as *Method 2* except using non-negative β_1 and β_2 .

We have the following observations.

1. *Method 2-g* is slower than *Method 2*. We have mentioned in the paper that the superiority of $\bar{\mathbf{d}}_k = \mathbf{d}_{k-1}$ may be because it comes from solving a sub-problem of using some second-order information. Further, $-\nabla f(\mathbf{w}_k)$, like \mathbf{d}_k , uses information of the current iteration while additional information from another subset of data is employed to find \mathbf{d}_{k-1} .

2. *Method 2-con* is slightly slower than *Method 2*. Although β_1 or β_2 by *Method 2* may be negative and cause some difficulties to interpret what the direction $\beta_1 \mathbf{d}_k + \beta_2 \bar{\mathbf{d}}_k$ is, they lead to the smallest second-order approximation of the function-value reduction. This property may explain *Method 2*'s faster convergence.

IV Details of Applying the Proposed Method to Maximum Entropy

For maximum entropy, the optimization problem is

$$\min_{\mathbf{w}} f(\mathbf{w}), \quad \text{where } f(\mathbf{w}) \equiv \frac{1}{2} \mathbf{w}^T \mathbf{w} + \frac{C}{l} \sum_{i=1}^l \left(\log \left(\sum_{c=1}^k \exp(\mathbf{w}_c^T \mathbf{x}_i) \right) - \mathbf{w}_{y_i}^T \mathbf{x}_i \right).$$

Let

$$P_{i,s} = \frac{\exp(\mathbf{w}_s^T \mathbf{x}_i)}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \mathbf{x}_i)}.$$

The gradient of $f(\mathbf{w})$ is

$$\nabla f(\mathbf{w}) = \begin{bmatrix} \nabla f^1(\mathbf{w}) \\ \vdots \\ \nabla f^k(\mathbf{w}) \end{bmatrix},$$

where

$$\nabla f^t(\mathbf{w}) = \mathbf{w}_t + \frac{C}{l} \left(\sum_{i=1}^l P_{i,t} \mathbf{x}_i - \sum_{i:y_i=t} \mathbf{x}_i \right) \in R^{n \times 1}.$$

Because $f(\mathbf{w})$ is twice continuously differentiable, the Hessian matrix of $f(\mathbf{w})$ can be derived.

- Case 1: When $t = s$,

$$\begin{aligned} \nabla^2 f(\mathbf{w})^{t,s} &= I_{n \times n} + \frac{C}{l} \sum_{i=1}^l \left(\frac{\exp(\mathbf{w}_t^T \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^T}{\sum_{c=1}^k \exp(\mathbf{w}_c^T \mathbf{x}_i)} - \frac{(\exp(\mathbf{w}_t^T \mathbf{x}_i) \mathbf{x}_i)(\exp(\mathbf{w}_s^T \mathbf{x}_i) \mathbf{x}_i)^T}{(\sum_{c=1}^k \exp(\mathbf{w}_c^T \mathbf{x}_i))^2} \right) \\ &= I_{n \times n} + \frac{C}{l} \sum_{i=1}^l \left((P_{i,t} \mathbf{x}_i) \mathbf{x}_i^T - (P_{i,t} \mathbf{x}_i)(P_{i,s} \mathbf{x}_i)^T \right) \in R^{n \times n}, \end{aligned}$$

where $I_{n \times n}$ is an identity matrix.

- Case 2: When $t \neq s$,

$$\begin{aligned} \nabla^2 f(\mathbf{w})^{t,s} &= \frac{C}{l} \sum_{i=1}^l \frac{-(\exp(\mathbf{w}_t^T \mathbf{x}_i) \mathbf{x}_i)(\exp(\mathbf{w}_s^T \mathbf{x}_i) \mathbf{x}_i)^T}{(\sum_{c=1}^k \exp(\mathbf{w}_c^T \mathbf{x}_i))^2} \\ &= \frac{C}{l} \sum_{i=1}^l (P_{i,t} \mathbf{x}_i)(P_{i,s} \mathbf{x}_i)^T \in R^{n \times n}. \end{aligned}$$

Therefore, the Hessian-vector product is

$$H\mathbf{v} = \begin{bmatrix} (H\mathbf{v})^1 \\ \vdots \\ (H\mathbf{v})^k \end{bmatrix}, \text{ where } \mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_k \end{bmatrix}$$

and

$$(H\mathbf{v})^t = \mathbf{v}_t + \frac{C}{l} \sum_{i=1}^l (P_{i,t}(\mathbf{v}_t^T \mathbf{x}_i - \sum_{c=1}^k P_{i,c} \mathbf{v}_c^T \mathbf{x}_i)) \mathbf{x}_i.$$

For the convergence analysis, we represent $\nabla^2 f(\mathbf{w})$ in a form similar to (29) for logistic regression.

$$\nabla^2 f(\mathbf{w}) = I_{kn \times kn} + \frac{C}{l} \bar{X}^T (D - E) \bar{X}, \quad (\text{IV.11})$$

where

$$\bar{X} = \begin{bmatrix} X & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & X & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & X \end{bmatrix} \in R^{kl \times kn},$$

$$D = \begin{bmatrix} D^1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & D^2 & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & D^k \end{bmatrix} \in R^{kl \times kl}, \text{ and } E = \begin{bmatrix} E^{11} & \cdots & \cdots & E^{1k} \\ E^{21} & E^{22} & \cdots & E^{2k} \\ \vdots & \vdots & \vdots & \vdots \\ E^{k1} & E^{k2} & \cdots & E^{kk} \end{bmatrix} \in R^{kl \times kl}.$$

(IV.12)

In (IV.12), D^s is a diagonal matrix with

$$D_{ii}^s = P_{i,s},$$

and E^{ts} is also a diagonal matrix with

$$E_{ii}^{ts} = P_{i,t} P_{i,s}.$$

Using (IV.11) we will prove Theorem 1 for the convergence. The only difference from that for logistic regression is on the boundedness of $\|H_{S_k}\|$ and $\|\nabla^2 f(\mathbf{w})\|$. From (IV.11) and the fact that $|D_{ii}^t| \leq 1$ and $|E_{ii}^{ts}| \leq 1, \forall i = 1, \dots, l, \|\nabla^2 f(\mathbf{w})\|$ is bounded by

$$\begin{aligned} \|\nabla^2 f(\mathbf{w})\| &\leq 1 + \frac{C}{l} \|\bar{X}^T\| \|D - E\| \|\bar{X}\| \\ &\leq 1 + \frac{C}{l} \|\bar{X}^T\| (\|D\| + \|-E\|) \|\bar{X}\| \\ &\leq 1 + \frac{C}{l} (1 + \sqrt{kl}) \|\bar{X}^T\| \|\bar{X}\|. \end{aligned}$$

For $\|H_{S_k}\|$, it is easy to have that

$$1 \leq \|H_{S_k}\| \leq \|\nabla^2 f(\mathbf{w}_k)\|.$$

V Details of Hessian-free Approaches for Neural Networks

The idea of calculating the Hessian-vector product is to define the following \mathcal{R} -operator for any function of $\boldsymbol{\theta}$ and then repeatedly apply the chain rule as follows:

$$\mathcal{R}_v\{f\} = \lim_{\epsilon \rightarrow 0} \frac{f(\boldsymbol{\theta} + \epsilon \mathbf{v}) - f(\boldsymbol{\theta})}{\epsilon}.$$

Then

$$H_k \mathbf{v} = \lim_{\epsilon \rightarrow 0} \frac{\nabla f(\boldsymbol{\theta} + \epsilon \mathbf{v}) - \nabla f(\boldsymbol{\theta})}{\epsilon} = \mathcal{R}_v\{\nabla f\}. \quad (\text{V.13})$$

For convenience, we replace \mathcal{R}_v with \mathcal{R} . To use (V.13), we first obtain ∇f by the following operations. Let \mathbf{x} and \mathbf{z} denote the \mathbf{x}^m and \mathbf{z}^m vectors of the m th layer, respectively. Further, we denote w_{ij} , $i = 1, \dots, n_m$, $j = 1, \dots, n_{m-1}$ as elements of the weight matrix W^m and let $\bar{\mathbf{z}}$ be the vector \mathbf{z}^{m-1} . Assume $\partial f / \partial z_i$, $i = 1, \dots, n_m$ are available. From (35), we have

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \frac{\partial f}{\partial z_i} \sigma'(x_i), \\ \frac{\partial f}{\partial w_{ij}} &= \frac{\partial f}{\partial x_i} \bar{z}_j, \\ \frac{\partial f}{\partial \bar{z}_j} &= \sum_{i=1}^{n_m} w_{ij} \frac{\partial f}{\partial x_i}. \end{aligned}$$

This backward process can be computed from the last to the first layer. In the end the collection of $\nabla_{W^1} f, \dots, \nabla_{W^L} f$ is $\nabla f(\boldsymbol{\theta})$.

To obtain $H_k \mathbf{v}$, we apply the \mathcal{R} -operator to the above terms and obtain

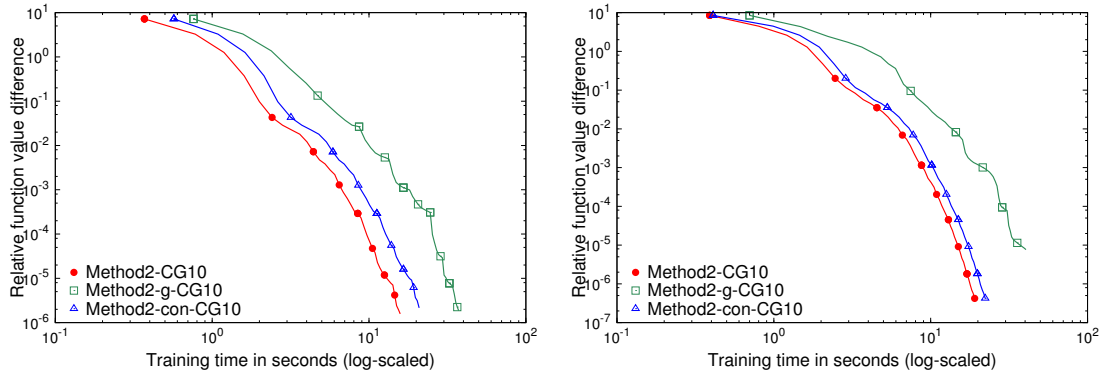
$$\begin{aligned} \mathcal{R}\left\{\frac{\partial f}{\partial x_i}\right\} &= \sigma'(x_i) \mathcal{R}\left\{\frac{\partial f}{\partial z_i}\right\} + \frac{\partial f}{\partial z_i} \sigma''(x_i) \mathcal{R}\{x_i\}, \\ \mathcal{R}\left\{\frac{\partial f}{\partial w_{ij}}\right\} &= \bar{z}_j \mathcal{R}\left\{\frac{\partial f}{\partial x_i}\right\} + \mathcal{R}\{\bar{z}_j\} \frac{\partial f}{\partial x_i}, \\ \mathcal{R}\left\{\frac{\partial f}{\partial \bar{z}_j}\right\} &= \sum_{i=1}^{n_m} (w_{ij} \mathcal{R}\left\{\frac{\partial f}{\partial x_i}\right\} + v_{ij} \frac{\partial f}{\partial x_i}), \end{aligned}$$

where v_{ij} is an element of the vector \mathbf{v} in (V.13). It corresponds to w_{ij} , so $\mathcal{R}(w_{ij}) = v_{ij}$. We see that $\mathcal{R}\{x_i\}$ and $\mathcal{R}\{\bar{z}_j\}$ are also needed. They are not computed in the backward process. Instead, we can pre-calculate them in the following forward process.

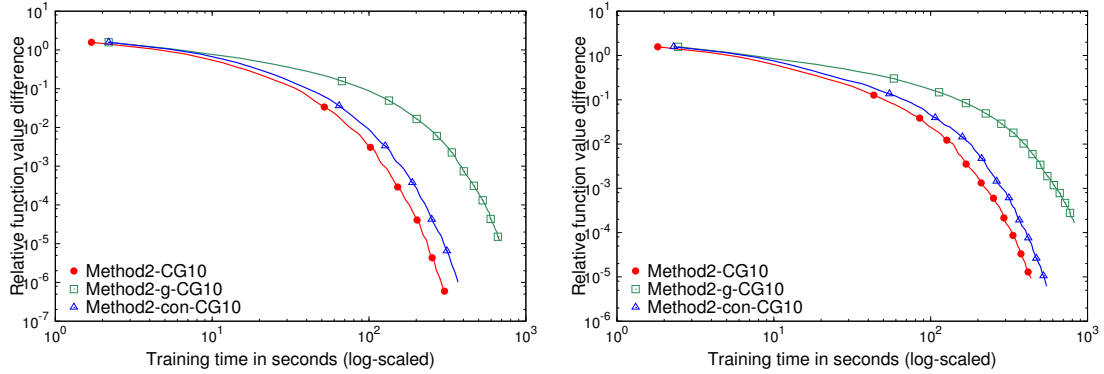
$$\begin{aligned} \mathcal{R}\{x_i\} &= \mathcal{R}\left\{\sum_{j=1}^{n_{m-1}} w_{ji} \bar{z}_j\right\} = \sum_{j=1}^{n_{m-1}} (w_{ji} \mathcal{R}\{\bar{z}_j\} + v_{ji} \bar{z}_j) \\ \mathcal{R}\{z_i\} &= \mathcal{R}\{\sigma(x_i)\} = \mathcal{R}\{x_i\} \sigma'(x_i). \end{aligned}$$

References

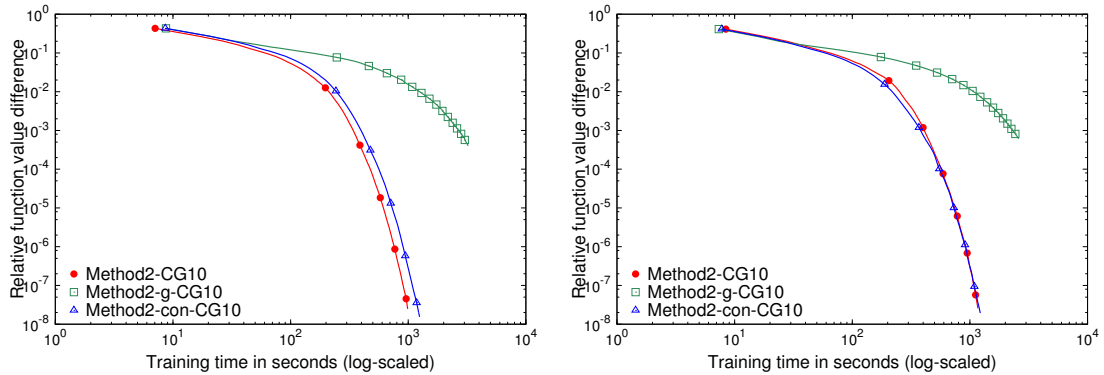
R. H. Byrd, G. M. Chin, W. Neveitt, and J. Nocedal. On the use of stochastic Hessian information in optimization methods for machine learning. *SIAM Journal on Optimization*, 21(3):977–995, 2011.



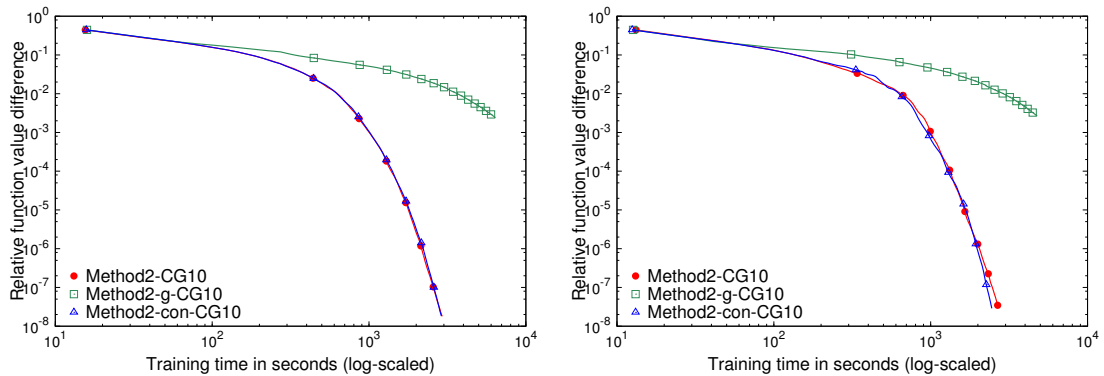
(a) news20



(b) yahoo-korea

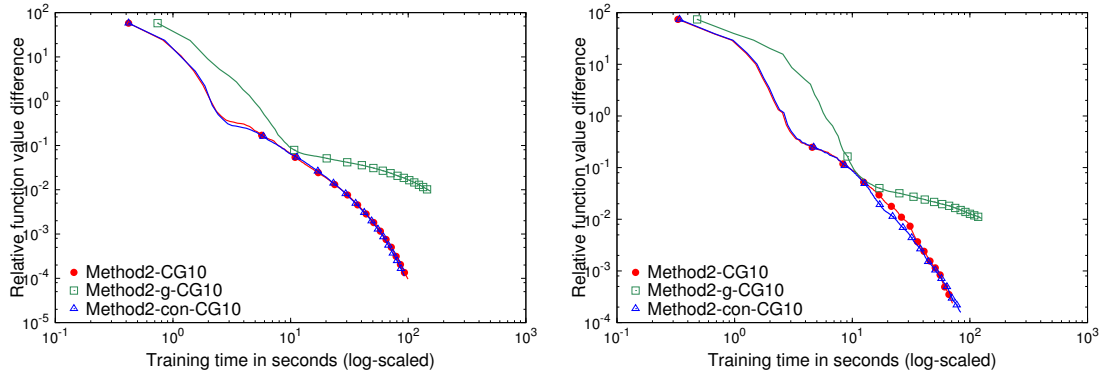


(c) kdd2010-a

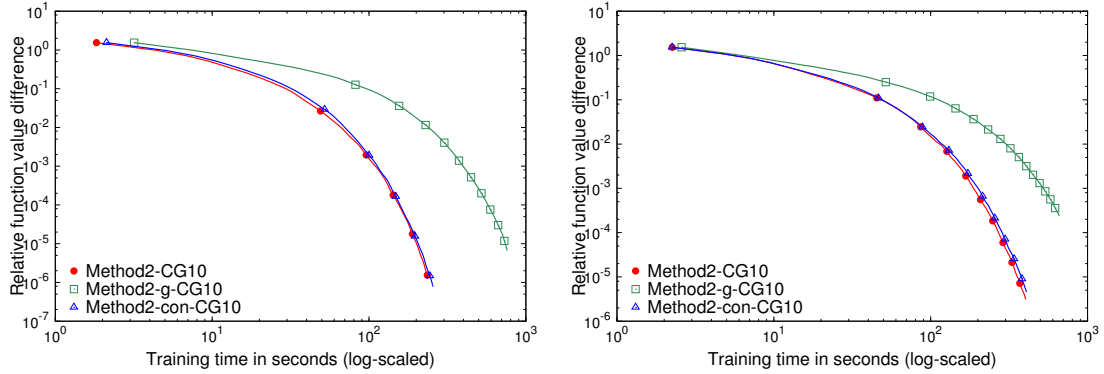


(d) kdd2010-b

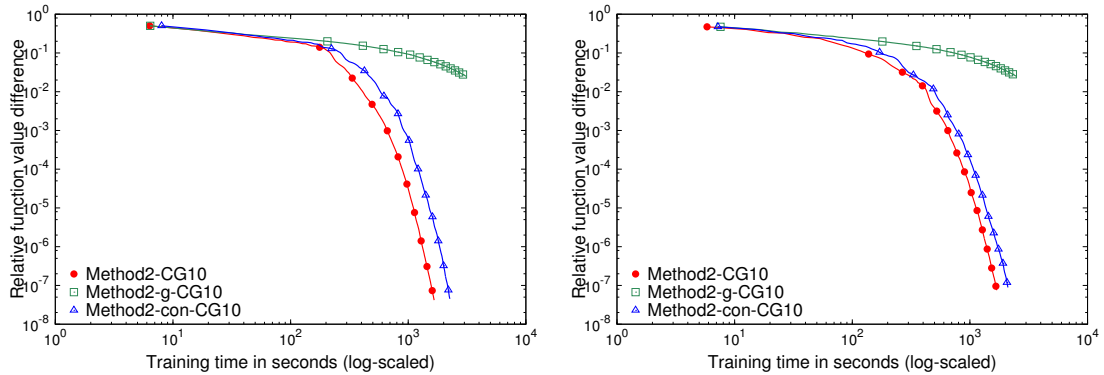
Figure III.2: Experiments on logistic regression using the three settings listed in Section III.3. We present running time (in seconds) versus the relative difference to the optimal function value. Both x -axis and y -axis are log-scaled. Left: $|S_k|/l = 5\%$. Right: $|S_k|/l = 1\%$.



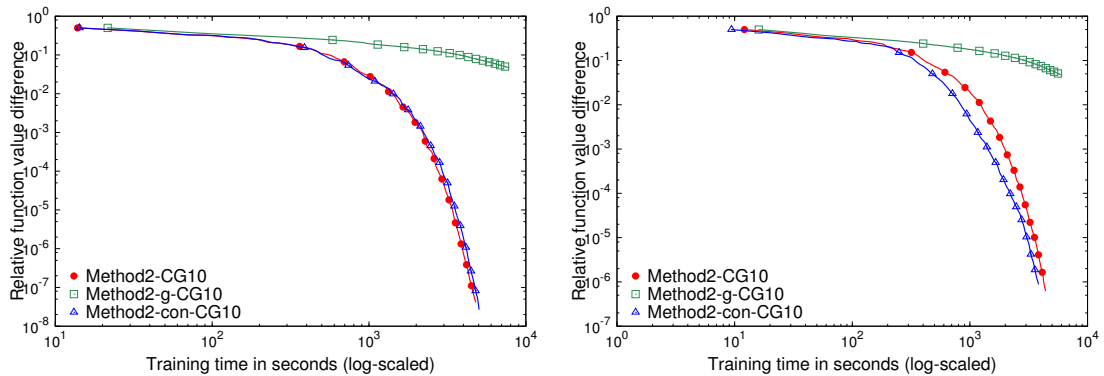
(a) news20



(b) yahoo-korea



(c) kdd2010-a



(d) kdd2010-b

Figure III.3: Experiments on l_2 -loss linear SVM using the three settings listed in Section III.3. We present running time (in seconds) versus the relative difference to the optimal function value. Both x -axis and y -axis are log-scaled. Left: $|S_k|/l = 5\%$. Right: $|S_k|/l = 1\%$.