

# An Efficient Alternating Newton Method for Learning Factorization Machines

WEI-SHENG CHIN, National Taiwan University  
 BO-WEN YUAN, National Taiwan University  
 MENG-YUAN YANG, National Taiwan University  
 CHIH-JEN LIN, National Taiwan University

To date, factorization machines (FM) have emerged as a powerful model in many applications. In this work, we study the training of FM with the logistic loss for binary classification, which is a non-linear extension of the linear model with the logistic loss (i.e., logistic regression). For the training of large-scale logistic regression, Newton methods have been shown to be an effective approach, but it is difficult to apply such methods to FM because of the non-convexity. We consider a modification of FM that is multi-block convex and propose an alternating minimization algorithm based on Newton methods. Some novel optimization techniques are introduced to reduce the running time. Our experiments demonstrate that the proposed algorithm is more efficient than stochastic gradient algorithms and coordinate descent methods. The parallelism of our method is also investigated for the acceleration in multi-threading environments.

CCS Concepts: •Computing methodologies → Factor analysis; Learning latent representations; Supervised learning by classification; •Mathematics of computing → Nonconvex optimization; •Information systems → Collaborative filtering; •Theory of computation → Shared memory algorithms;

Additional Key Words and Phrases: Newton methods, preconditioned conjugate gradient methods, sub-sampled Hessian matrix

## ACM Reference Format:

Wei-Sheng Chin, Bo-Wen Yuan, Meng-Yuan Yang, and Chih-Jen Lin, 2016. An Efficient Alternating Newton Method for Learning Factorization Machines. *ACM Trans. Intell. Syst. Technol.* 0, 0, Article 0 (0), 32 pages. DOI: 0000001.0000001

## 1. INTRODUCTION

Binary classification has been an important technique for many practical applications. Assume that the observations and the scalar result of an event are respectively described by an  $n$ -dimensional feature vector  $x$  and a label  $y \in \{-1, 1\}$ . The model is an output function from the feature space to a real number,  $\hat{y} : \mathcal{R}^n \rightarrow \mathcal{R}$ , where the output label is the sign of the output value. Among many existing models, we are interested in factorization machines (FM), which have been shown to be useful for high dimensional and sparse data sets [Rendle 2010]. Assume that we have  $l$  training events,

$$(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l),$$

This work is supported by MOST of Taiwan via grants 104-2622-E-002-012-CC2 and 104-2221-E-002-047-MY3 and MOE of Taiwan via grants 104R7872 and 105R7872.

Authors' email: {d01944006,r03944049,b01902037,cjlin}@csie.ntu.edu.tw

Author's addresses: Wei-Sheng Chin, Graduate Institute of Networking and Multimedia, National Taiwan University; Bo-Wen Yuan, Meng-Yuan Yang, and Chih-Jen Lin, Department of Computer Science, National Taiwan University

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 0 ACM. 2157-6904/0/-ART0 \$15.00

DOI: 0000001.0000001

where  $y_i$  and  $x_i$  are the label and the feature vector of the  $i$ th training instance, respectively. For an instance  $(y, x)$ , the output value of FM is defined by

$$\hat{y}(x) = w^T x + \sum_{j=1}^n \sum_{j'=j+1}^n u_j^T u_{j'} x_j x_{j'}, \quad (1)$$

where  $w \in \mathcal{R}^n$  and  $U = [u_1, \dots, u_n] \in \mathcal{R}^{d \times n}$  are the parameters of FM. The  $j$ th column of  $U$ ,  $u_j$ , can be viewed as the  $j$ th feature's representation in an latent space with the dimensionality  $d$  specified by the user. In (1), the variable  $w$  describes the linear relation between the input features and the output label, while the second term of (1) includes all feature conjunctions that capture the information carried by feature co-occurrences (also called feature interactions or feature conjunctions). Moreover, the coefficient of the conjunction between  $x_j$  and  $x_{j'}$  is determined by the inner product of  $u_j$  and  $u_{j'}$ . To determine FM's parameters, we solve the following optimization problem.

$$\min_{w, U} \quad \frac{\lambda}{2} \|w\|^2 + \frac{\lambda'}{2} \|U\|_F^2 + \sum_{i=1}^l \xi(\hat{y}(x_i); y_i). \quad (2)$$

FM's objective function consists of the squared sum of all model parameters for preventing over-fitting and the sum of  $l$  loss values. The loss function  $\xi(\cdot)$  encourages the consistency between  $\hat{y}(x)$  and the actual label  $y$ . For practical classification problems, one can use, for example, the logistic loss or the squared loss:

$$\xi_{\log}(\hat{y}; y) = \log(1 + e^{-y\hat{y}}) \quad (3)$$

$$\xi_{\text{sq}}(\hat{y}; y) = \frac{1}{2}(\hat{y} - y)^2. \quad (4)$$

Note that the regression-oriented loss (4) can be adopted here because it panellizes the distance between the predicted value and the true label and therefore ensure the consistency between them. Furthermore, if  $U$  is not considered, using any of (3) and (4) leads to logistic regression (or LR for short) and linear regression, respectively.

By following the concept of computing the conjunction coefficient of a feature pair from their latent representations, some FM variants have been proposed. For example, Blondel et al. [2016] proposed modifying the output function by introducing some new variables,  $v_1, \dots, v_n$ , and change output value in (1) to

$$\begin{aligned} \hat{y}(x) &= w^T x + \frac{1}{2} \sum_{j=1}^n \sum_{j'=1}^n u_j^T v_{j'} x_j x_{j'} \\ &= w^T x + \frac{1}{2} \underbrace{(Ux)^T (Vx)}_{x^T U^T V x}, \end{aligned} \quad (5)$$

where  $V = [v_1, \dots, v_n] \in \mathcal{R}^{d \times n}$  encodes the feature representations in another latent space. They replace  $U^T V$  with  $\frac{1}{2}(U^T V + V^T U)$  in (5), but both of the output value and the optimization problem remain unchanged. In addition to the appearance of another latent space associated with  $V$ , (5) differs from (1) in the following aspects. First, the multiplications  $x_j x_{j'}$  and  $x_{j'} x_j$  are treated as different feature conjunctions, so they are characterized by different coefficients. For example, to calculate the inner product in determining the coefficient of a conjunction  $x_j x_{j'}$  in (5), the latent representation of the left feature  $u_j$  is picked up from the original latent space and that of the right feature  $v_{j'}$  is in the new latent space. Second, a feature index can occur twice in one conjunction; that is, self-interaction is allowed. Resulting from these two changes, the

Table I: FM solvers which have been studied.

Solver \ Loss	FM type	Squared loss	Logistic loss
SG	(2)	[Rendle 2010]	[Rendle 2010; Ta 2015; Juan et al. 2016]
	(6)	[Blondel et al. 2016]	
MCMC	(2)	[Rendle 2012]	
CD	(2)	[Rendle 2012]	[Blondel et al. 2016]
	(6)	[Blondel et al. 2016]	
L-BFGS	(2), (6)	[Blondel et al. 2016]	

summation in (5) must iterate through all indexes in  $\{(j, j') \mid j, j' = 1, \dots, n\}$ . Besides, Blondel et al. [2016] do not consider the  $1/2$  coefficient of the non-linear term in (5). We added it because the new formulation nearly doubles the number of feature conjunctions in (1). Finally, the new optimization problem associated with (5) is clearly formed by three blocks of variables  $\mathbf{w}$ ,  $U$ , and  $V$ :

$$\min_{\mathbf{w}, U, V} F(\mathbf{w}, U, V), \quad (6)$$

where

$$F(\mathbf{w}, U, V) = \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{\lambda'}{2} \|U\|_F^2 + \frac{\lambda''}{2} \|V\|_F^2 + \sum_{i=1}^l \xi \left( y_i \left( \mathbf{w}^T \mathbf{x}_i + \frac{1}{2} (U \mathbf{x}_i)^T (V \mathbf{x}_i) \right) \right). \quad (7)$$

In this paper, we mainly focus on (5) because some of its nice properties needed in developing our method.

Problems (2) and (5) are not easy to solve because of their non-convex objective functions. When (4) is used, stochastic gradient methods (SG) and coordinate descent methods (CD) are two of the most effective ones [Rendle 2012; Blondel et al. 2016]. Also, Blondel et al. [2016] demonstrated that a classical quasi-Newton method, L-BFGS, is largely outperformed by CD and SG. For training FM with the logistic loss, most packages use stochastic methods. For example, a Markov Chain Monte Carlo method (MCMC) and a classical stochastic gradient method in LIBFM [Rendle 2012] while LIBFFM [Juan et al. 2016] implements ADAGRAD, an advanced SG by Duchi et al. [2011]. Ta [2015] also reported the effectiveness of a variant of ADAGRAD when considering (3). See Table I for existing FM algorithms and the loss functions they have been implemented for. On the other hand, [Blondel et al. 2016] argued that their CD procedure can be adopted to deal with the logistic loss with some minor modifications, but neither implementation nor experimental result is available to confirm it. Because (3) is very popular in several key applications (e.g., click-through rate prediction), we think it's important to examine if the state of the art deterministic algorithm (i.e., CD) is applicable or SG just dominates here. Unfortunately, the situation may not be very optimistic for CD. It has been known that to LR CD becomes less attractive because of needing a relatively large number of exponential/logarithmic operations comparing with standard arithmetic operations [Yuan et al. 2012a]. On the other hand, Yuan et al. [2012a] show that Newton-type methods requiring fewer exponential/logarithmic operations are compelling. Thus, we believe that the possibility of designing a Newton-based method to beat CD when considering (3) still exists. Note that L-BFGS can induce less exponential/logarithmic operations, but we do not include it in this study because of the following two reasons. First, L-BFGS is not an ideal solver candidate as it is much worse than both of SG and CD in learning FM with the squared loss. Second,

Lin et al. [2008] demonstrated that their truncated Newton method can surpassed L-BFGS in training LR [Lin et al. 2008]. Our contributions are summarized below.

- We explore using Newton method in training FM in detail. A Newton-based algorithm is proposed along with some techniques for a large reduction on its running time.
- The parallelism of our method is also investigated for the acceleration in multi-threading environments.
- We empirically demonstrate the extraordinary efficiency of the proposed algorithm against stochastic gradient methods and coordinate descent methods, targeting (6) with (3).

We emphasize that our discussion on existing and the proposed FM solvers is applicable to both of (4) and (3). Nevertheless, we merely consider (3) in our experiments because we are motivated by the difficulty of training FM when the logistic loss is used.

The paper is organized as follows. Section 2 reviews existing state of the art FM solvers. We develop an efficient algorithm based on Newton methods in Section 3. Section 4 discusses some algorithmic issues and implementation techniques for our method. Then, we study its parallelization of our method in Section 5. The experimental results in Section 6 confirm the effectiveness of the proposed method. Finally, Section 7 concludes our work.

## 2. STATE OF THE ART METHODS FOR TRAINING FACTORIZATION MACHINES

In this section, we discuss the two state of the art FM solvers in detail, CD and SG. Our focus is (6) with the loss functions defined in (3) and (4). The reason of not considering (2) is that some of (6)'s nice properties are required in developing our method.

### 2.1. Stochastic Gradient Methods

It is mentioned in Section 1 that stochastic gradient methods are very popular in training FM. Among existing solvers, we adopt ADAGRAD because its superiority over LIBFM has been established in [Juan et al. 2016].

Consider the aforementioned  $l$  training instances  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, l$  again. ADAGRAD solves optimization problems in the following form.

$$\min_{\tilde{\mathbf{w}} \in \mathcal{R}^{\tilde{n}}} \sum_{i=1}^l f(\tilde{\mathbf{w}}; y_i, \mathbf{x}_i), \quad (8)$$

where  $\tilde{\mathbf{w}}$  is the unknown variable and  $f(\tilde{\mathbf{w}}; y_i, \mathbf{x}_i)$  is the cost function that the  $i$ th instance incurs. At the  $k$ th iteration, an instance  $(y_i, \mathbf{x}_i)$  is sampled from the  $l$  training instances to calculate the stochastic gradient,  $\tilde{\mathbf{g}}^k = \nabla f(\tilde{\mathbf{w}}^k; y_i, \mathbf{x}_i)$ . Then, the current solution is updated via

$$\tilde{\mathbf{w}}^{k+1} \leftarrow \tilde{\mathbf{w}}^k - \eta_0 \tilde{G}^{-\frac{1}{2}} \tilde{\mathbf{g}}^k,$$

where  $\eta_0 > 0$  is a pre-specified constant and  $\tilde{G}^{-\frac{1}{2}}$  is the inverse of the root of a diagonal matrix defined by

$$\tilde{G} = \sum_{k'=1}^k \text{diag} \left( \tilde{\mathbf{g}}^{k'} \right)^2.$$

The complete procedure that ADAGRAD solves (8) is shown in Algorithm 1.

**Algorithm 1** ADAGRAD.

---

```

1: Given initial solution  $\tilde{w}^0$ , number of iterations  $\bar{k}$ , and  $\eta_0 > 0$ .
2: Initialization  $\tilde{G}$  with zeros.
3: for  $k \leftarrow 1, \dots, \bar{k}$  do
4:   Draw an index  $i \in \{1, \dots, l\}$ 
5:    $\tilde{g} \leftarrow \nabla f(\tilde{w}^k; y_i, \mathbf{x}_i)$ 
6:    $\tilde{G} \leftarrow \tilde{G} + \text{diag}(\tilde{g})^2$ 
7:    $\tilde{w}^{k+1} \leftarrow \tilde{w}^k - \eta_0 \tilde{G}^{-1/2} \tilde{g}$ 
8: end for

```

---

To apply ADAGRAD to learn the considered FM, we begin with rewriting (6) into a form of (8). Let

$$\tilde{w} = [\mathbf{w}^T, \mathbf{u}_1^T, \dots, \mathbf{u}_n^T, \mathbf{v}_1^T, \dots, \mathbf{v}_n^T]^T \in \mathcal{R}^{n+2dn}. \quad (9)$$

We argue that if

$$f(\tilde{w}; y_i, \mathbf{x}_i) = \sum_{j \in N_i} \left( \frac{\lambda}{2|\Omega_j|} w_j^2 + \frac{\lambda'}{2|\Omega_j|} \|\mathbf{u}_j\|^2 + \frac{\lambda''}{2|\Omega_j|} \|\mathbf{v}_j\|^2 \right) + \xi(\hat{y}(\mathbf{x}_i); y_i), \quad (10)$$

then (8) is equivalent to (6), where  $N_i$  is the index set of the  $i$ th instance's non-zero features,  $\Omega_j$  is the indexes set of the instances whose  $j$ th features are not zero, and  $|\cdot|$  returns the size when the input is a set. Because the summation of the loss terms in (10) over  $i = 1, \dots, l$  is obviously the loss term in (7), we only check if the regularization terms in (10) lead to the same regularization terms in (7). First, for  $w$ 's regularization, we have

$$\sum_{i=1}^l \sum_{j \in N_i} \frac{\lambda}{2|\Omega_j|} w_j^2 = \frac{\lambda}{2} \sum_{j=1}^n \sum_{i \in \Omega_j} \frac{1}{|\Omega_j|} w_j^2 = \frac{\lambda}{2} \sum_{j=1}^n w_j^2 = \frac{\lambda}{2} \|\mathbf{w}\|^2.$$

By an analogous derivation, inversely scaling  $\|\mathbf{u}_j\|^2$  and  $\|\mathbf{v}_j\|^2$  by  $|\Omega_j|$  produces the desired regularization in (7). The equivalence between (10) and (7) is therefore verified.

From (10) and the definition of  $\tilde{w}$  in (9), we deduce the stochastic gradient with respect to the model parameters in each iteration. For  $w$ , we have

$$\frac{\partial f(\tilde{w}; y_i, \mathbf{x}_i)}{\partial w_j} = \begin{cases} \frac{\lambda}{|\Omega_j|} w_j + \xi'(\hat{y}(\mathbf{x}_i); y_i) x_{ij} & \text{if } j \in N_i, \\ 0 & \text{otherwise,} \end{cases} \quad (11)$$

where  $\xi'(\hat{y}; y)$  is the derivative of the selected loss function with respect to  $\hat{y}$ . Note that

$$\begin{aligned} \xi'_{\log}(\hat{y}; y) &= \frac{-y}{1 + e^{y\hat{y}}} \\ \xi'_{\text{sq}}(\hat{y}; y) &= \hat{y} - y. \end{aligned} \quad (12)$$

Taking a closer look at (12), one may quickly realizes that computing the derivative of (3) is much more slower than calculating the other components. The chief reason is that exponential and logarithmic function evaluations are much more expensive than standard arithmetic operations. In contrast, (4)'s evaluation is not especially high-priced. To simplify subsequent notations, we define  $\mathbf{p}_i = U\mathbf{x}_i$  and  $\mathbf{q}_i = V\mathbf{x}_i$ . Taking the partial derivative of (10) with respect to  $\mathbf{u}_j$  and  $\mathbf{v}_j$  yields the derivatives along other

coordinates.

$$\begin{aligned}\frac{\partial f(\tilde{\mathbf{w}}; y_i, \mathbf{x}_i)}{\partial \mathbf{u}_j} &= \begin{cases} \frac{\lambda'}{|\Omega_j|} \mathbf{u}_j + \frac{1}{2} \xi'(\hat{y}(\mathbf{x}_i); y_i) \mathbf{q}_i x_{ij} & \text{if } j \in N_i, \\ \mathbf{0}_{d \times 1} & \text{otherwise.} \end{cases} \\ \frac{\partial f(\tilde{\mathbf{w}}; y_i, \mathbf{x}_i)}{\partial \mathbf{v}_j} &= \begin{cases} \frac{\lambda''}{|\Omega_j|} \mathbf{v}_j + \frac{1}{2} \xi'(\hat{y}(\mathbf{x}_i); y_i) \mathbf{p}_i x_{ij} & \text{if } j \in N_i, \\ \mathbf{0}_{d \times 1} & \text{otherwise.} \end{cases}\end{aligned}\quad (13)$$

Substituting (11) and (13) into the step that calculates the stochastic gradient in Algorithm 1, we get the procedure to solve (6) by ADAGRAD.

Next, we analyze required computational complexity in one for-loop iteration in Algorithm 1. Provided that the cost of a loss-related function call (i.e., the evaluation of  $\xi(\hat{y}; y)$  or its derivatives) can be much more higher than a standard arithmetic operation in some circumstances, the amounts of loss-related function calls and other operations are separately discussed. If the instance  $(y_i, \mathbf{x}_i)$  is drawn, we first spend  $\mathcal{O}(|N_i|d)$  operations to calculate and store  $\mathbf{p}_i, \mathbf{q}_i$  and then the output value  $\hat{y}_i$  via

$$(5) = \mathbf{w}^T \mathbf{x}_i + \frac{1}{2} \mathbf{p}_i^T \mathbf{q}_i. \quad (14)$$

The following step is substituting  $\hat{y}_i$  into  $\xi'(\hat{y}_i; y)$  and caching its output. Next, we use (11)-(13) to compute the  $(|N_i| + 2|N_i|d)$  non-zero values of  $\tilde{\mathbf{g}}$ . Then, the diagonal elements of  $\tilde{\mathbf{G}}$  correspond to the non-zero coordinates in  $\tilde{\mathbf{g}}$  are updated. Finally, new values are computed and assigned to the coordinates of  $\tilde{\mathbf{w}}$  linked to the non-zero coordinates in  $\tilde{\mathbf{g}}$ . The total arithmetic complexity in one iteration is therefore  $\mathcal{O}(|N_i|d)$  with only one call to the loss function's derivative. Let  $(\# \text{ nnz})$  denote the total number of non-zeros in  $X$ . To process one epoch (i.e.,  $l$  instances), we expect  $\mathcal{O}(d \times (\# \text{ nnz}))$  operations and  $l$  loss-related function calls.

From a practical perspective, selecting proper storage formats for the used variables is important in obtaining an efficient implementation. A well-known principle is to use one leading to shorter physical distance between two consecutively-accessed variables. In Algorithm (3), the major variables are  $X, U$ , and  $V$ . For each of them, we can use one of row-major or column-major formats. Since ADAGRAD accesses  $X$  instance-by-instance, it should be stored in a row-major format (e.g., compressed row format because  $X$  is usually sparse in large-scale problems). Because the coordinates in  $\mathbf{u}_j, j \in N_i$  are sequentially accessed when computing (13), a column-major format is suggested for  $U$ . Similarly,  $V$  is stored in column-major format. It follows that the matrices of accumulated gradients are also in column-major.

## 2.2. Coordinate Descent Methods

Coordinate descent methods, a family of iterative procedures, have gained great successes in solving large-scale minimization problems. At each iteration, only one variable is updated by solving an associated sub-problem while other variables are fixed. For the update sequence, here we choose a cyclic one following [Blondel et al. 2016]. That is, we cyclically go through

$$w_1, \dots, w_j, \underbrace{U_{11}, U_{12}, \dots, U_{1n}}_{\text{the 1st row of } U}, \underbrace{V_{11}, V_{12}, \dots, V_{1n}}_{\text{the 1st row of } V}, \dots, \underbrace{U_{d1}, U_{d2}, \dots, U_{dn}}_{\text{the } d\text{th row of } U}, \underbrace{V_{d1}, V_{d2}, \dots, V_{dn}}_{\text{the } d\text{th row of } V}. \quad (15)$$

For forming a sub-problem and obtaining its solution, we leverage the concept in state of the art methods developed for LR and support vector machine [Chang et al. 2008; Huang et al. 2010]. The situations for updating the coordinates of  $\mathbf{w}, U$ , and  $V$  are quite similar. Let us first outline the update of  $U_{cj}$ , the value at the  $c$ th row and the

$j$ th column of  $U$ , and then extend the result to other cases. By fixing all variables but  $U_{cj}$  in (6), we obtain  $U_{cj}$ 's sub-problem,

$$\min_{U_{cj} \in \mathcal{R}} F_{cj}^U(U_{cj}), \quad (16)$$

where

$$F_{cj}^U(U_{cj}) = \frac{\lambda'}{2} U_{cj}^2 + \sum_{i \in \Omega_j} \xi(y_i; \hat{y}_i).$$

For solving (16), a single-variable Newton method is employed; that is, if the current solution is  $U_{cj}$ , the update direction would be

$$s = -g/h, \quad (17)$$

where

$$\begin{aligned} g &= \frac{\partial F_{cj}^U}{\partial U_{cj}} = \lambda' U_{cj} + \frac{1}{2} \sum_{i \in \Omega_j} \xi'(\hat{y}_i; y_i) (\mathbf{q}_i)_c (\mathbf{x}_i)_j \\ h &= \frac{\partial^2 F_{cj}^U}{\partial U_{cj}^2} = \lambda' + \frac{1}{4} \sum_{i \in \Omega_j} \xi''(\hat{y}_i; y_i) (\mathbf{q}_i)_c^2 (\mathbf{x}_i)_j^2. \end{aligned} \quad (18)$$

For the second-order derivatives, we note

$$\begin{aligned} \xi''_{\log}(\hat{y}; y) &= \frac{e^{y\hat{y}}}{(1 + e^{y\hat{y}})^2} \\ \xi''_{\text{sq}}(\hat{y}; y) &= 1. \end{aligned} \quad (19)$$

If  $P$  and  $Q$  are available,  $\mathcal{O}(|\Omega_j|)$  arithmetic and loss-related operations are involved in the evaluation of (18). Then, to guarantee a sufficient objective function reduction and hence the convergence, a back-tracking line search is launched to find the largest step size  $\theta \in \{1, \beta, \beta^2, \dots\}$  satisfying Armijo rule,

$$F_{cj}^U(U_{cj} + \theta s) - F_{cj}^U(U_{cj}) \leq \theta \nu g s, \quad (20)$$

where  $\nu \in (0, 1)$  and  $\beta \in (0, 1)$  are pre-specified constants. If  $\hat{y}_i, i \in \Omega_j$  are available, computing the left-hand side in (20) can be finished within  $\mathcal{O}(|\Omega_j|)$  arithmetic operations and loss-related function calls. Precisely,

$$\begin{aligned} F_{cj}^U(U_{cj} + \theta s) - F_{cj}^U(U_{cj}) &= \frac{\lambda'}{2} (2\theta U_{cj} s + \theta^2 s^2) + \sum_{i \in \Omega_j} \xi(\hat{y}_i + \theta \Delta_i; y_i) \\ &\quad - \sum_{i \in \Omega_j} \xi(\hat{y}_i; y_i), \end{aligned} \quad (21)$$

where  $\Delta_i = 0.5 \times s(\mathbf{q}_i)_c (\mathbf{x}_i)_j$  because in (5), the relevant output part is

$$(\mathbf{U}\mathbf{x})^T (\mathbf{V}\mathbf{x}) = \mathbf{q}^T \mathbf{U}\mathbf{x} = \sum_{c=1}^d \sum_{j=1}^n (\mathbf{q})_c U_{cj} (\mathbf{x})_j$$

and only  $U_{cj}$  is being adjusted along  $s$ . In addition,  $\Delta_i, i \in \Omega_j$  can be pre-computed in the beginning and reused in checking every  $\theta$ . Once  $\theta$  is found, the variable  $U_{cj}$  is adjusted via

$$U_{cj} \leftarrow U_{cj} + \theta s.$$

**Algorithm 2** CD for solving (16).

---

```

1: Given iteration number  $\bar{k}_{\text{inner}}$ ,  $\lambda' > 0$ ,  $0 < \nu < 1$ ,  $0 < \beta < 1$ , the current  $U_{cj}$ ,  $\hat{y}_i$ ,  $(\mathbf{p}_i)_c$ ,
   and  $(\mathbf{q}_i)_c$ ,  $i = 1, \dots, l$ .
2: for  $k \leftarrow \{1, \dots, \bar{k}_{\text{inner}}\}$  do
3:   Evaluate (17) to obtain  $s$ 
4:    $\Delta_i \leftarrow s(\mathbf{q}_i)_c(\mathbf{x}_i)_j/2$ ,  $\forall i \in \Omega_j$ 
5:   for  $\theta \leftarrow \{1, \beta, \beta^2, \dots\}$  do
6:     if (21)  $\leq \theta \nu g s$  then
7:        $U_{cj} \leftarrow U_{cj} + \theta s$ 
8:       Maintain variables via (22),  $\forall i \in \Omega_j$ .
9:       break
10:    end if
11:  end for
12: end for
13: Output  $U_{cj}$  as the solution of (16).
```

---

It is known that for an efficient CD implementation,  $\mathbf{p}_i$ ,  $\mathbf{q}_i$ , and  $\hat{y}_i$ ,  $\forall i$  should be cached and carefully maintained through iterations to save computation [Blondel et al. 2016]. The update of  $U_{cj}$  causes that the  $c$ th coordinate at  $\mathbf{p}_i$  and  $\hat{y}_i$ ,  $i \in \Omega_j$  are changed. Recomputing  $\mathbf{p}_i$  and  $\hat{y}_i$  for each  $i \in \Omega_j$  may need  $\mathcal{O}(|N_i|d)$  operations, but all necessary changes can be done via a  $\mathcal{O}(|\Omega_j|)$  rule below.

$$\begin{aligned} (\mathbf{p}_i)_c &\leftarrow (\mathbf{p}_i)_c + \theta s(\mathbf{x}_i)_j, \quad \forall i \in \Omega_j. \\ \hat{y}_i &\leftarrow \hat{y}_i + \theta \Delta_i \end{aligned} \quad (22)$$

Algorithm 2 sketches our CD procedure of solving (16). Besides, if  $\xi(\hat{y}; y)$  is  $\mu$ -smooth, letting

$$h = \lambda' + \frac{\mu}{4} \sum_{i \in \Omega_j} (\mathbf{q}_i)_c^2 (\mathbf{x}_i)_j^2 \quad (23)$$

can guarantee the monotonic decrease of the objective function with  $\theta = 1$ . It leads to the CD procedure without line search in [Blondel et al. 2016]. Nevertheless, replacing the calculation of  $h$  in (18) with (23) does not change the algorithm's computation complexity, so we use (18) for more accurate second-order information following a state of the art LR solver [Huang et al. 2010]. All iterative algorithms need a stopping condition. For Algorithm 2, a fixed number of iterations  $\bar{k}_{\text{inner}}$  is set. Concerning (18), (21), and (22), Algorithm 2 yields a computational cost at

$$\mathcal{O}(\bar{k}_{\text{inner}} \times (\# \text{ of line search iterations} + 2) \times |\Omega_j|), \quad (24)$$

where the coefficient “2” accounts for the evaluation of (18) and (22). If all coordinates of  $U$  are updated once, the induced complexity is

$$\mathcal{O}(d\bar{k}_{\text{inner}} \times (\# \text{ of expected line search iterations} + 2) \times (\# \text{ nnz})). \quad (25)$$

As none of loss-related operations is present in (22), the total number of loss-related function calls in Algorithm 2 and updating  $U$  once can be simply obtained by changing the coefficient “2” in (24) and (25) to “1,” respectively.

With some minor changes, Algorithm 2 can be reused to update the other coordinates. For  $V_{cj}$ , we just need to swap  $(\lambda', U_{cj}, (\mathbf{p}_i)_c, (\mathbf{q}_i)_c)$  and  $(\lambda'', V_{cj}, (\mathbf{q}_i)_c, (\mathbf{p}_i)_c)$ . For  $w_j$ , the strategy is to replace  $(\lambda', U_{cj}, (\mathbf{p}_i)_c, (\mathbf{q}_i)_c)$  with  $(\lambda, w_j, 2, 2)$  and disable the adjustment of  $(\mathbf{p}_i)_c$  in (22). All combined, Algorithm 3 demonstrates the final two-layer CD framework for solving (6). Because the adjustment of  $U$  and  $V$  costs more than the



**Algorithm 3** Solving (6) via CD.

---

```

1: Given an initial solution  $(w, U, V)$  and number of outer iterations  $\bar{k}_{\text{outer}}$ .
2: Calculate and cache  $p_i, q_i$ , and  $\hat{y}_i = w^T x_i + \frac{1}{2} p_i^T q_i, i = 1, \dots, l$ .
3: for  $k \leftarrow \{1, \dots, \bar{k}_{\text{outer}}\}$  do
4:   for  $j \leftarrow \{1, \dots, n\}$  do
5:     Update  $w_j$  via Algorithm 2 by replacing  $(\lambda', U_{cj}, (q_i)_c)$  with  $(\lambda, w_j, 2)$  and discarding the first rule in (22).
6:   end for
7:   for  $c \leftarrow \{1, \dots, d\}$  do
8:     for  $j \leftarrow \{1, \dots, n\}$  do
9:       Update  $U_{cj}$  via Algorithm 2.
10:    end for
11:    for  $j \leftarrow \{1, \dots, n\}$  do
12:      Update  $V_{cj}$  via Algorithm 2 by swapping  $(\lambda', U_{cj}, (p_i)_c)$  and  $(\lambda'', V_{cj}, (q_i)_c)$ .
13:    end for
14:  end for
15: end for

```

---

update of  $w$ , the complexity of updating all coordinates once (called an outer iteration) is (25) with loss-related function calls at the same level. Besides, in Algorithm (3), the initialization step requires  $\mathcal{O}(d \times (\# \text{nnz}))$  extra operations to prepare  $p_i, q_i$ , and  $\hat{y}_i, i = 1, \dots, l$ .

Finally, we discuss which of row- and column-major orders should be obeyed when storing the elements in  $X, U, V, P$ , and  $Q$ . Recall the three major steps in Algorithm 2: (18), the line search procedure (including the preparation of  $\Delta_i, \forall i \in \Omega_j$ ), and (22). For each step, we observe that aligning elements in  $Q$  in row-major order may increase memory continuity when accessing  $(q_i)_c$ , resulting from its summations over the instance indexes in  $\Omega_j$ . This argument also holds for  $P$  when updating  $V$ . Thus, both of  $P$  and  $Q$  should be stored in row-major format. Because we are always looping through  $(x_i)_j, i \in \Omega_j, X$  should be in column-major format. From the update sequence specified by (15), row-major format is recommended for both of  $U$  and  $V$ .

### 3. ALTERNATING NEWTON METHODS FOR LEARNING FACTORIZATION MACHINES

From Sections 2.1-2.2, we know that CD's and SG's ratios of loss-related function calls to standard arithmetic operations are  $\mathcal{O}(1)$  and  $\mathcal{O}\left(\frac{l}{d \times \# \text{nnz}}\right)$ , respectively. It implies that when the loss function is composed of some expensive operations, the running time of CD can be heavily prolonged. For Newton-type methods, solving the direction-finding sub-problem is usually the heaviest step in each iteration and, fortunately, costs no loss-related operation. Thus, it may be a more suitable method for tackling (3) than CD. However, the non-convexity of (2)-6 brings some difficulties in developing a Newton-type FM solver. The Hessian matrix (the second-order derivative) is not always positive definite, so a direct application of the Newton method can fail to find a descent direction. We may consider a positive-definite approximation of the Hessian matrix such as the Gauss-Newton matrix [Schraudolph 2002] used in deep learning [Martens 2010], but such an approach is generally less effective than Newton methods for convex optimization problems. On the other hand, for problems related to matrix factorization, the objective function is multi-block convex, so alternating Newton methods have been widely applied. That is, sequentially a convex sub-problem of one block of variables is solved by Newton method while other variables are fixed. In fact, if the squared loss is used, then the alternating Newton method is reduced to the alter-

---

**Algorithm 4** Solving the modified FM problem (6) via alternating minimization.

---

```

1: Given an initial solution  $(\mathbf{w}, U, V)$ 
2: while stopping condition is not satisfied do
3:    $\mathbf{w} \leftarrow \arg \min_{\mathbf{w}} F(\mathbf{w}, U, V)$ 
4:    $U \leftarrow \arg \min_U F(\mathbf{w}, U, V)$ 
5:    $V \leftarrow \arg \min_V F(\mathbf{w}, U, V)$ 
6: end while

```

---

nating least square method (e.g., [Zhou et al. 2008]) because Newton method solves a quadratic problem in one iteration.

Unfortunately, the original FM problem (2) is not in a multi-block convex form so that each block contains many variables. The reason is that every  $\mathbf{u}_j^T \mathbf{u}_{j'}$  in (1) is a non-convex term with respect to  $\mathbf{u}_j$  and  $\mathbf{u}_{j'}$ . Then, we can only get a small convex sub-problem by fixing all but one  $\mathbf{u}_j$ . On the other hand, it has been mentioned in [Blondel et al. 2016] that (7) is a multi-block convex function of  $\mathbf{w}$ ,  $U$ , and  $V$ . That is, (7) becomes a convex sub-problem if two of  $\mathbf{w}$ ,  $U$ , and  $V$  are fixed. Then the idea of alternating minimization by using any convex optimization algorithm to solve each sub-problem can be applied. In the rest of this paper we consider the optimization problem (6) of using the new model in (5) rather than the original FM optimization problem.

For the update sequence of the block minimization, it is possible to have a dynamic scheme, but for the sake of simplicity we focus on the cyclic one. Iteratively the following three optimization sub-problems are solved.

$$\min_{\mathbf{w}} F(\mathbf{w}, U, V), \quad (26)$$

$$\min_U F(\mathbf{w}, U, V), \quad (27)$$

$$\min_V F(\mathbf{w}, U, V). \quad (28)$$

The overall procedure to minimize (6) is summarized in Algorithm 4. In the rest of this section, we discuss how to apply Newton-type methods to solve the sub-problems. The resulting procedure is called ANT (alternating Newton method). We begin with showing that each sub-problem is in a form similar to linear classification problem.

We emphasize that our method is especially designed to alleviate CD's weakness induced by using (3), through the subsequently discussion is applicable to (4) as well.

### 3.1. Relation Between (6)'s Sub-problems and Linear Classification

Although (26)-(28) are different optimization problems, we show that they are all linear classification problems by the change of variables. Given  $l$  training instances, a linear classification problem can be written as

$$\min_{\tilde{\mathbf{w}} \in \mathcal{R}^{\tilde{n}}} \tilde{f}(\tilde{\mathbf{w}}), \quad (29)$$

where

$$\begin{aligned} \tilde{f}(\tilde{\mathbf{w}}) &= \frac{\tilde{\lambda}}{2} \|\tilde{\mathbf{w}}\|^2 + \sum_{i=1}^l \xi(\tilde{y}_i; y_i), \\ \tilde{y}_i &= \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i + \tilde{c}_i, \end{aligned} \quad (30)$$

$\tilde{\mathbf{x}}_i$  is the feature vector of the  $i$ th instance, and  $\tilde{c}_i$  is a constant.<sup>1</sup>

Before showing some reformulations of (26)-(28), recall (14). It is straightforward to see that (26) is in a form of (29) when

$$\tilde{\lambda} = \lambda, \quad \tilde{\mathbf{w}} = \mathbf{w}, \quad \tilde{\mathbf{x}}_i = \mathbf{x}_i, \quad \text{and} \quad \tilde{c}_i = \frac{1}{2} \mathbf{p}_i^T \mathbf{q}_i.$$

Note that constants such as  $\|U\|_F^2$  and  $\|V\|_F^2$  are not considered when updating  $\mathbf{w}$ . To write (27) in a form of (29), we note that when  $\mathbf{w}$  and  $V$  are fixed, (14) can be reformulated as an affine function of  $U$ ,

$$\tilde{y}_i = \frac{1}{2} (U \mathbf{x}_i)^T \mathbf{q}_i + \mathbf{w}^T \mathbf{x}_i = \text{vec}(U)^T \left( \frac{1}{2} (\mathbf{x}_i \otimes \mathbf{q}_i) \right) + \mathbf{w}^T \mathbf{x}_i,$$

where the vectorization operator outputs a column vector by stacking the input matrix's columns and " $\otimes$ " denotes the Kronecker product. With  $\|\text{vec}(U)\|^2 = \|U\|_F^2$ , (27) is in the form of (29) by

$$\tilde{\lambda} = \lambda', \quad \tilde{\mathbf{w}} = \text{vec}(U), \quad \tilde{\mathbf{x}}_i = \frac{1}{2} (\mathbf{x}_i \otimes \mathbf{q}_i), \quad \text{and} \quad \tilde{c}_i = \mathbf{w}^T \mathbf{x}_i. \quad (31)$$

Similarly, (29) becomes (28) if

$$\tilde{\lambda} = \lambda'', \quad \tilde{\mathbf{w}} = \text{vec}(V), \quad \tilde{\mathbf{x}}_i = \frac{1}{2} (\mathbf{x}_i \otimes \mathbf{p}_i), \quad \text{and} \quad \tilde{c}_i = \mathbf{w}^T \mathbf{x}_i. \quad (32)$$

We can conclude that (29) can be used to represent all sub-problems in Algorithm 4, so the same optimization algorithm can be used to solve them. We give the gradient and the Hessian-matrix of (29) below for applying optimization algorithms.

$$\nabla \tilde{f}(\tilde{\mathbf{w}}) = \tilde{\lambda} \tilde{\mathbf{w}} + \tilde{X}^T \mathbf{b} \quad (33)$$

$$\nabla^2 \tilde{f}(\tilde{\mathbf{w}}) = \tilde{\lambda} I + \tilde{X}^T D \tilde{X}, \quad (34)$$

where  $\tilde{X} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_l]^T \in \mathcal{R}^{l \times \tilde{n}}$  is the data matrix,  $\mathbf{b} = [\xi'(\tilde{y}_1; y_1), \dots, \xi'(\tilde{y}_l; y_l)]^T \in \mathcal{R}^l$ , and  $D$  is a diagonal matrix whose  $i$ th diagonal element is  $\xi''(\tilde{y}_i; y_i)$ . See (12) and (19) for  $\xi'(\tilde{y}; y)$ 's and  $\xi''(\tilde{y}; y)$ 's definitions. With  $\tilde{\lambda} > 0$ , it is clear that (34) is positive definite, so the optimal solution of (29) is unique.

### 3.2. Truncated Newton Methods for Solving (29)

We discuss a truncated Newton method to solve (29). At the  $k$ th iteration, Newton method obtains an update direction by minimizing the second-order Taylor expansion at  $\tilde{f}(\tilde{\mathbf{w}}^k)$ .

$$\min_{\tilde{\mathbf{s}} \in \mathcal{R}^{\tilde{n}}} \tilde{\mathbf{g}}^T \tilde{\mathbf{s}} + \frac{1}{2} \tilde{\mathbf{s}}^T \tilde{H} \tilde{\mathbf{s}}, \quad (35)$$

where  $\tilde{\mathbf{w}}^k$  is the solution to be improved,  $\tilde{H} = \nabla^2 \tilde{f}(\tilde{\mathbf{w}}^k)$ , and  $\tilde{\mathbf{g}} = \nabla \tilde{f}(\tilde{\mathbf{w}}^k)$ . Because  $\tilde{H}$  is positive definite, (35) is equivalent to a linear system,

$$\tilde{H} \tilde{\mathbf{s}} = -\tilde{\mathbf{g}}. \quad (36)$$

In practice, (36) may not need to be exactly solved, especially in early iterations. Therefore, truncated Newton methods have been introduced to approximately solve (36) while maintain the convergence [Nash 2000]. Iterative methods such as conjugate gradient method (CG) are often used for approximately solving (36), so at each iteration

<sup>1</sup>In standard linear classification problems,  $\tilde{c}_i = 0$ . Note that we do not consider a bias term in the output function.

an inner iterative process is involved. Another difficulty to solve (36) is that storing the Hessian matrix  $\tilde{H}$  in  $\mathcal{O}(\tilde{n}^2)$  space is not possible if  $\tilde{n}$  is large. To overcome the memory issue, Komarek and Moore [2005] and Keerthi and DeCoste [2005] showed that for problems like (29), the conjugate gradient method that mainly requires a sequence of Hessian-vector products can be performed without explicitly forming the Hessian matrix:

$$\tilde{H}\tilde{s} = \lambda\tilde{s} + \tilde{X}^T(D(\tilde{X}\tilde{s})), \quad (37)$$

where  $\tilde{s}$  is the iterate in the CG procedure. In [Lin et al. 2008], this idea is further incorporated into a truncated Newton framework.

Like coordinate descent methods and many other optimization algorithms, after a direction is found in our Newton method, we must decide the step size taken along that direction. Here we consider the back-tracking line-search procedure mentioned in Section 2.2 and discuss a trick sharing the same spirit with (21) for efficient computation. Concerning (29) and the  $\tilde{n}$ -dimensional search direction, the line search's stopping condition specified in (20) becomes

$$\tilde{f}(\tilde{w}^k + \theta\tilde{s}) - \tilde{f}(\tilde{w}^k) \leq \theta\nu\tilde{g}^T\tilde{s}. \quad (38)$$

Recalculating the function value at each  $\tilde{w}^k + \theta\tilde{s}$  is expensive, where the main cost is on calculating  $(\tilde{w} + \theta\tilde{s})^T\tilde{x}_i, \forall i$ . However, for linear models, the following trick in [Yuan et al. 2010] can be employed. Assume that

$$\tilde{y}_i = (\tilde{w}^k)^T\tilde{x}_i + \tilde{c}_i \text{ and } \tilde{\Delta}_i = \tilde{s}^T\tilde{x}_i, \quad i = 1, \dots, l \quad (39)$$

are available. At an arbitrary  $\theta$ , we can calculate

$$(\tilde{w}^k + \theta\tilde{s})^T\tilde{x}_i + \tilde{c}_i = \tilde{y}_i + \theta\tilde{\Delta}_i \quad (40)$$

to get the new output value. By further maintaining

$$(\tilde{w}^k)^T\tilde{s}, \quad \|\tilde{s}\|^2, \text{ and } \tilde{g}^T\tilde{s}, \quad (41)$$

the total cost of checking the condition of (38) is merely  $\mathcal{O}(l)$  because

$$\begin{aligned} \tilde{f}(\tilde{w}^k + \theta\tilde{s}) - \tilde{f}(\tilde{w}^k) &= \frac{\tilde{\lambda}}{2} \left( 2\theta(\tilde{w}^k)^T\tilde{s} + \theta^2\|\tilde{s}\|^2 \right) + \sum_{i=1}^l \xi(\tilde{y}_i + \theta\tilde{\Delta}_i; y_i) \\ &\quad - \sum_{i=1}^l \xi(\tilde{y}_i; y_i). \end{aligned} \quad (42)$$

In the end of the line-search procedure, the model is updated by

$$\tilde{w}^{k+1} \leftarrow \tilde{w}^k + \theta\tilde{s},$$

and the last value obtained in (40) can be used as the new  $\tilde{y}_i$  in the next iteration; see (39). When  $\tilde{y}_i, \forall i$  are available, we immediately obtain the elements in  $b$  and  $D$  according to (12) and (19), respectively. Notice that if the logistic loss is considered, a temporal variable  $e^{y_i\tilde{y}_i}$  which participates in the computation of  $b_i$  can be cached and reused in the computation of  $D_{ii}$  to save some exponential operations.

Algorithm 5 summarizes the above procedure. The computational complexity of one iteration in Algorithm 5 can be estimated by

$$(\text{cost of } \tilde{f} \text{ and } \tilde{g}) + (\# \text{ of CG iterations}) \times (\text{cost of } \tilde{H}\tilde{s}) + (\# \text{ of line search iterations}) \times \mathcal{O}(l). \quad (43)$$

---

**Algorithm 5** A CG-based truncated Newton method with line search for solving (29).

---

```

1: Given initial value of  $\tilde{w}^0$  and  $0 < \tilde{\epsilon} < 1$ .
2: Compute and cache  $\tilde{y}_i, \forall i$ .
3:  $f^0 \leftarrow \tilde{f}(\tilde{w}^0)$ 
4: for  $k \leftarrow \{0, 1, \dots\}$  do
5:   Calculate and store  $b_i$  (and  $e^{y_i \tilde{y}_i}$  if (3) is the loss) using (12),  $\forall i$ .
6:    $\tilde{g} \leftarrow \tilde{\lambda} \tilde{w} + \tilde{X}^T \mathbf{b}$ 
7:   if  $k = 0$  then
8:      $\|\tilde{g}^0\| \leftarrow \|\tilde{g}\|$ 
9:   end if
10:  if  $\|\tilde{g}\| \leq \tilde{\epsilon} \|\tilde{g}^0\|$  then
11:    Output  $\tilde{w}^k$  as the solution.
12:    break
13:  end if
14:  Compute  $D_{ii}$  via (19),  $\forall i$ .
15:  Solve (36) approximately via CG to get an update direction  $\tilde{s}$ .
16:  Calculate variables listed in (41) and  $\tilde{\Delta}_i, \forall i$ .
17:  for  $\theta \leftarrow \{1, \beta, \beta^2, \dots\}$  do
18:    Compute  $\delta = f(\tilde{w}^k + \theta \tilde{s}) - f^k$  via (42).
19:    if  $\delta \leq \nu \theta \tilde{g}^T \tilde{s}$  then
20:       $\tilde{w}^{k+1} \leftarrow \tilde{w}^k + \theta \tilde{s}$ 
21:       $f^{k+1} \leftarrow f^k + \delta$ 
22:      Let the last value calculated in (40) be the next  $\tilde{y}_i, \forall i$ 
23:      break
24:    end if
25:  end for
26: end for

```

---

Note that the above line-search procedure is also responsible for calculating the next function value, so in (43), except the  $\mathcal{O}(l)$  cost for calculating (40) at each line-search step, we consider all initialization tasks such as calculating  $\tilde{\Delta}_i, \forall i$  in (39) and (41) as the cost of evaluating the objective function. The complexity of evaluating the objective function also absorbs the cost of calculating  $b_i$  and  $D_{ii}, i = 1, \dots, l$ . Since there are  $\mathcal{O}(l)$  at each of Lines 5, 14, and 18, the number of total loss-related function calls in one iteration is

$$\mathcal{O}((\# \text{ of line search iterations} + 2) \times l). \quad (44)$$

Although we do not exactly solve (36), Algorithm 5 asymptotically converges to the optimal point if  $\tilde{s}$  satisfies

$$\|\tilde{H}\tilde{s} + \tilde{g}\| \leq \eta \|\tilde{g}\|,$$

where  $0 < \eta < 1$  can be either a pre-specified constant or a number controlled by a dynamic strategy [Eisenstat and Walker 1994; Nash and Sofer 1996].

Based on our discussion, Algorithm 5 is able to solve the three sub-problems because they are all in the form of (29). However, without taking their problem structures into account, the implementation may be inefficient. In particular, (27) and (28) significantly differ from traditional linear classification problems because they are matrix-variable problems. In Section 3.3, we will discuss some techniques for an efficient implementation of the truncated Newton method.

### 3.3. An Efficient Implementation of Algorithm 5 for Solving Sub-problems (27)-(28)

Among the three sub-problems, (26) is very close to standard linear classification problem, so the implementation is similar to past works. However, (27) and (28) are more different because their variables are matrices. Here, we only discuss details of solving (27) because the situation for (28) is similar.

To begin with, we check the function evaluation, in which the main task is on calculating  $\tilde{y}_i, \forall i$ . From (30) and (31),

$$\begin{aligned}\tilde{y}_i &= \text{vec}(U)^T \tilde{\mathbf{x}}_i + \tilde{c}_i \\ &= \frac{1}{2} \text{vec}(U)^T (\mathbf{x}_i \otimes \mathbf{q}_i) + \tilde{c}_i \\ &= \frac{1}{2} \mathbf{q}_i^T U \mathbf{x}_i + \tilde{c}_i, \quad i = 1, \dots, l,\end{aligned}$$

or equivalently

$$\begin{aligned}\tilde{\mathbf{y}} &= \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_l \end{bmatrix} = \tilde{X} \text{vec}(U) + \tilde{\mathbf{c}} \\ &= \frac{1}{2} (Q^T .* (XU^T)) \mathbf{1}_{d \times 1} + \tilde{\mathbf{c}},\end{aligned}\tag{45}$$

where  $Q = [\mathbf{q}_1, \dots, \mathbf{q}_l] \in \mathcal{R}^{d \times l}$ , “ $.*$ ” stands for the element-wise product of two matrices, and  $\mathbf{1}_{d \times 1}$  is the vector of ones. For (45), we can calculate and store  $XU^T$  as a dense matrix regardless of whether  $X$  is sparse or not; the cost is  $\mathcal{O}(d \times (\# \text{ nnz}))$ . Similar to (45),  $\tilde{\Delta}$  used in line search can be computed via

$$\tilde{\Delta}_i = \frac{1}{2} \mathbf{q}_i^T S \mathbf{x}_i, \quad \forall i \text{ or equivalently } \tilde{\Delta} = \tilde{X} \tilde{\mathbf{s}} = \frac{1}{2} (Q^T .* (XS^T)) \mathbf{1}_{d \times 1},\tag{46}$$

where  $\tilde{\mathbf{s}} = \text{vec}(S)$  with  $S \in \mathcal{R}^{d \times n}$  is the search direction.

Then, we investigate how to compute the gradient of (27). By some reformulations, we have

$$\begin{aligned}\nabla \tilde{f}(\text{vec}(U)) &= \lambda' \text{vec}(U) + \tilde{X}^T \mathbf{b} \\ &= \lambda' \text{vec}(U) + \frac{1}{2} \sum_{i=1}^l (b_i \mathbf{x}_i \otimes \mathbf{q}_i) \\ &= \lambda' \text{vec}(U) + \frac{1}{2} \sum_{i=1}^l \text{vec}(b_i \mathbf{q}_i \mathbf{x}_i^T) \\ &= \lambda' \text{vec}(U) + \frac{1}{2} \text{vec} \left( \sum_{i=1}^l (\mathbf{q}_i b_i \mathbf{x}_i^T) \right) \\ &= \lambda' \text{vec}(U) + \frac{1}{2} \text{vec}(Q \text{diag}(\mathbf{b})X),\end{aligned}\tag{47}$$

where  $\text{diag}(\mathbf{b})$  is a diagonal matrix in which the  $i$ th diagonal element is  $b_i$ . Depending on the size of  $Q$  and  $X$ , or whether  $X$  is sparse, we can decide to calculate  $Q \text{diag}(\mathbf{b})$  or  $\text{diag}(\mathbf{b})X$  first. Then, the main operation for gradient evaluation is a matrix-matrix product that costs  $\mathcal{O}(d \times (\# \text{ nnz}))$ .

Next we discuss Hessian-vector products in CG. From (31), (37), and  $\tilde{s} = \text{vec}(S)$ ,

$$\begin{aligned}\tilde{H}\tilde{s} &= \lambda' \text{vec}(S) + \tilde{X}^T(D(\tilde{X} \text{vec}(S))) \\ &= \lambda' \text{vec}(S) + \tilde{X}^T(Dz)\end{aligned}\quad (48)$$

$$\begin{aligned}&= \lambda' \text{vec}(S) + \sum_{i=1}^l (Dz)_i \tilde{x}_i \\ &= \lambda' \text{vec}(S) + \sum_{i=1}^l (Dz)_i (\mathbf{x}_i \otimes \mathbf{q}_i) \\ &= \lambda' \text{vec}(S) + \frac{1}{2} \sum_{i=1}^l \text{vec}(\mathbf{q}_i (Dz)_i \mathbf{x}_i^T)\end{aligned}\quad (49)$$

$$= \lambda' \text{vec}(S) + \frac{1}{2} \text{vec}(Q \text{diag}(Dz)X), \quad (50)$$

where

$$z_i = \tilde{x}_i^T \tilde{s} = \frac{1}{2} (\mathbf{x}_i^T \otimes \mathbf{q}_i^T) \text{vec}(S) = \frac{1}{2} \mathbf{q}_i^T S \mathbf{x}_i, \quad i = 1, \dots, l \quad (51)$$

or equivalently

$$\mathbf{z} = \tilde{X} \text{vec}(S) = \frac{1}{2} (Q^T .* (XS^T)) \mathbf{1}_{d \times 1} \quad (52)$$

by following the same calculation in (46). Note that  $Dz$  is the product between a diagonal matrix  $D$  and a vector  $\mathbf{z}$ . We investigate the complexity of the Hessian-vector product. The first term in (48) is a standard vector scaling that costs  $\mathcal{O}(nd)$ . For the second term in (48), we separately consider two major steps:

- (1) The computation of  $\mathbf{z} = \tilde{X} \text{vec}(S)$ .
- (2) The product between  $\tilde{X}^T$  and a dense vector  $Dz$ .

In (52), we have seen that the first step can be done by the same way to compute (46). For the second step, we can refer to the calculation of  $\tilde{X}^T \mathbf{b}$  in (47). Therefore, following earlier discussion, the cost of one Hessian-vector product is

$$\mathcal{O}(d \times (\# \text{nnz})). \quad (53)$$

A disadvantage of the above setting is that the data matrix  $X$  is accessed twice in the two different steps. From (49) and (51), the Hessian-vector product can be represented as

$$\lambda' \text{vec}(S) + \frac{1}{4} \sum_{i=1}^l \text{vec}(\mathbf{q}_i D_{ii} (\mathbf{q}_i^T S \mathbf{x}_i) \mathbf{x}_i^T). \quad (54)$$

From this instance-wise representation, the second term can be computed by a single loop over  $\mathbf{x}_1, \dots, \mathbf{x}_l$ . While  $\mathbf{x}_i$  is still used twice:  $S \mathbf{x}_i$  and  $\mathbf{q}_i^T S \mathbf{x}_i$ , between them no other  $\mathbf{x}_j, j \neq i$  are accessed and therefore  $\mathbf{x}_i$  tends to remain at a higher layer of the memory hierarchy. Such a setting of getting better data locality has been considered for training LR [Zhuang et al. 2015]. However, because of the instance-wise setting, we can not benefit from optimized dense matrix operations (e.g., optimized BLAS) when computing, for example,  $SX^T$  in (52). Therefore, the decision of using (54) or not may hinge on whether  $X$  is a sparse matrix.

In line search, besides the calculation of  $\tilde{\Delta}$  in (46), from (41) we also need

$$\langle U, S \rangle, \|S\|_F^2, \text{ and } \langle G, S \rangle, \quad (55)$$

where  $G$  is the matrix form of (27)'s gradient and  $\langle \cdot, \cdot \rangle$  is the inner product of two matrices.

Algorithm 6 summarizes the details in our Newton procedure for solving (27). Matrix representations are used, so we never have to reshape  $S$  to  $\text{vec}(S)$  or  $U$  to  $\text{vec}(U)$ . The latest  $U$  in ANT is used as the initial point, so we can continue to improve upon the current solution. In Algorithm 6, the computational complexity per Newton iteration is

$$(\# \text{ of CG iterations} + 2) \times \mathcal{O}(d \times (\# \text{ nnz})) + (\# \text{ of line search iterations}) \times \mathcal{O}(l), \quad (56)$$

in which the term  $2 \times \mathcal{O}(d \times (\# \text{ nnz}))$  is from the costs of (46) and (47) for function and gradient evaluation. Note that the  $\mathcal{O}(nd)$  cost for the quantities in (55) and other places is not listed here because it is no more than  $\mathcal{O}(d \times (\# \text{ nnz}))$ . Let  $\bar{n} = \frac{\# \text{ nnz}}{l}$  denote the average non-zero features per instance. By (44) and (56), there are around

$$\frac{(\# \text{ of CG iteration} + 2) \times d\bar{n} + (\# \text{ of line search iterations})}{(\# \text{ of line search iterations} + 2)}$$

standard read, write, or arithmetic operations per loss-related function call. If the number of line search iterations is not larger than the number of CG iterations, the ratio of loss-related operations to other standard operations in ANT would be much smaller than  $\mathcal{O}(1)$  so that the major drawback in CD is alleviated considerably.

### 3.4. Convergence Guarantee and Stopping Conditions

Algorithm 8 is under the framework of block coordinate descent method (BCD), which sequentially update one block of variables while keeping all remaining blocks unchanged. A worthy-remained example of BCD is CD where each block contains only a single variable. See [Bertsekas 1999] a comprehensive introduction about BCD. It is known that if each sub-problem of a block has a unique optimal solution and is exactly solved, then the procedure converges to a stationary point [Bertsekas 1999, Proposition 2.7.1]. Each of our sub-problems (26)-(28) possesses a unique optimal solution because of the strongly convex regularization term. Therefore, if we exactly solve every sub-problem, then our method is guaranteed to converge to a stationary point of (6).

In practice, any optimization method needs a stopping condition. Here, we consider a relative setting

$$\|\nabla F(\mathbf{w}, U, V)\| \leq \epsilon \|\nabla F(\mathbf{w}_{\text{init}}, U_{\text{init}}, V_{\text{init}})\|, \quad (57)$$

where  $0 < \epsilon < 1$  is a small stopping tolerance and  $(\mathbf{w}_{\text{init}}, U_{\text{init}}, V_{\text{init}})$  indicates the initial point of the model. The use of this type of criteria includes, for example, [Lin 2007] for non-negative matrix factorization.

Each sub-problem in the alternating minimization procedure requires a stopping condition as well. A careful design is needed as otherwise the global stopping condition in (57) may never be reached. In Section 3.2, we consider a relative stopping condition,

$$\|\nabla \tilde{f}(\tilde{\mathbf{w}})\| \leq \tilde{\epsilon} \|\nabla \tilde{f}(\tilde{\mathbf{w}}^0)\|, \quad (58)$$

where  $\tilde{\mathbf{w}}^0$  is the initial point in solving the sub-problem. Take the sub-problem of  $U$  as an example. The gradient norm of (6) with respect to  $U$  is identical to the gradient norm of  $U$ 's sub-problem:

$$\|\nabla_U f(\mathbf{w}, U, V)\|_F = \|\nabla \tilde{f}(\text{vec}(U))\|.$$



---

**Algorithm 6** An implementation of Algorithm 5 for solving (27) by operations on matrix variables without vectorizing them.

---

```

1: Given  $0 < \tilde{\epsilon} < 1$  and the current  $(w, U, V)$ .
2:  $Q \leftarrow VX^T$ 
3: Compute and cache  $\tilde{y} = \frac{1}{2} (Q^T .* (XU^T)) \mathbf{1}_{d \times 1} + Xw$ .
4:  $f \leftarrow \frac{\lambda'}{2} \|U\|_F^2 + \sum_{i=1}^l \xi(\tilde{y}_i; y_i)$ .
5: for  $k \leftarrow \{0, 1, \dots\}$  do
6:   Calculate and store  $\tilde{y}_i$  (and  $e^{y_i \tilde{y}_i}$  if (3) is used) and then obtain  $b_i$  by (12),  $\forall i$ .
7:    $G \leftarrow \lambda' U + \frac{1}{2} Q \text{diag}(\mathbf{b}) X$ 
8:   if  $k = 0$  then
9:      $\|G^0\|_F \leftarrow \|G\|_F$ 
10:  end if
11:  if  $\|G\|_F \leq \tilde{\epsilon} \|G^0\|_F$  then
12:    Output  $U$  as the solution of (27).
13:    break
14:  end if
15:  Compute  $D_{ii}$  via (19),  $\forall i$ .
16:  Solve (36) approximately via CG to get an update direction  $S$ .
17:  Prepare variables listed in (55) and  $\tilde{\Delta} = \frac{1}{2} (Q^T .* (XS^T)) \mathbf{1}_{d \times 1}$ 
18:  for  $\theta \leftarrow \{1, \beta, \beta^2, \dots\}$  do
19:     $\delta \leftarrow \frac{\lambda'}{2} (2\theta \langle U, S \rangle + \theta^2 \|S\|_F^2) + \sum_{i=1}^l \xi(\tilde{y}_i + \theta \tilde{\Delta}_i; y_i) - \sum_{i=1}^l \xi(\tilde{y}_i; y_i)$ 
20:    if  $\delta \leq \theta \nu \langle G, S \rangle$  then
21:       $U \leftarrow U + \theta S$ 
22:       $f \leftarrow f + \delta$ 
23:       $\tilde{y} \leftarrow \tilde{y} + \theta \tilde{\Delta}$ 
24:      break
25:    end if
26:  end for
27: end for

```

---

If  $\|\nabla_U f(w, U, V)\|_F > 0$ , Algorithm 6 conducts at least one Newton step and the function value of  $F(w, U, V)$  must be decreased. The reason is that (58) is violated if we do not change  $U$ . In other words, any block with a non-vanished gradient would be adjusted. Therefore, unless

$$\|\nabla_w f(w, U, V)\| = \|\nabla_U f(w, U, V)\|_F = \|\nabla_V f(w, U, V)\|_F = 0,$$

one of the three blocks must be updated. That is, if the condition in (57) is not satisfied yet, our alternating minimization procedure must continue to update  $(w, U, V)$  rather than stay at the same point.

#### 4. TECHNIQUES TO ACCELERATE TRUNCATED NEWTON METHODS FOR SOLVING THE SUB-PROBLEMS

From (56), the computational bottleneck in our truncated Newton method is the CG procedure, so we consider two popular techniques, preconditioning and sub-sampled Hessian, for its acceleration.

##### 4.1. Preconditioned Conjugate Gradient Methods

The first technique aims to reduce the number of CG iterations for solving (36) by considering an equivalent but better conditioned linear system. Precisely, we consider a preconditioner  $\tilde{M}$  to approximately factorize  $\tilde{H}$  such that  $\tilde{H} \approx \tilde{M}\tilde{M}^T$ , and then use

CG to solve

$$\hat{H}\hat{s} = \hat{g}, \quad (59)$$

where  $\hat{H} = \tilde{M}^{-1}\tilde{H}\tilde{M}^{-T}$  and  $\hat{g} = \tilde{M}^{-1}\tilde{g}$ . Once  $\hat{s}$  is found, the solution of (36) can be recovered by  $\tilde{s} = \tilde{M}^{-T}\hat{s}$ . If  $\tilde{H}$  is perfectly factorized (i.e.,  $\tilde{H} = \tilde{M}\tilde{M}^T$ ), (59) can be solved in no time because  $\hat{H}$  is an identity matrix.

Many preconditioners have been proposed, e.g., diagonal preconditioner, incomplete Cholesky factorization, and polynomial preconditioners [Golub and Van Loan 2012]. However, finding a suitable preconditioner is not easy because first each CG iteration becomes more expensive and second the decrease of the number of CG iterations is not theoretically guaranteed. For a CG-based truncated Newton method for LR in [Lin et al. 2008], the use of preconditioned conjugate gradient methods (PCG) has been studied. They point out that the implicit use of  $\tilde{H}$  further increases the difficulty to implement a preconditioner and finally choose the simple diagonal preconditioner for the following two reasons. First, the diagonal elements of  $\tilde{H}$  can be constructed cheaply. Second, the extra computation introduced by the diagonal preconditioner in each CG iteration is relatively small. Experiments in [Lin et al. 2008] show that for LR, using diagonal preconditioners may not be always effective in decreasing the total number of CG iterations and the running time. Nevertheless, we think it is worth trying preconditioned CG here because of some differences between their settings and ours. In the alternating minimization procedure, we solve a sequence of optimization sub-problems of (6) rather than a single linear classification problem in [Lin et al. 2008]. In the early stage of the ANT, we only loosely solve the sub-problems and the effect of PCG may vary depending on the strictness of the stopping condition.

To see the details of using PCG in ANT, we consider the following diagonal preconditioner for solving (36) as an example.

$$\tilde{M} = \tilde{M}^T = \text{diag}(\tilde{h}), \quad (60)$$

where

$$\tilde{h} = \sqrt{\tilde{\lambda}\mathbf{1}_{\tilde{n} \times 1} + \sum_{i=1}^l D_{ii}(\tilde{\mathbf{x}}_i * \tilde{\mathbf{x}}_i)}.$$

Note that “ $\sqrt{\cdot}$ ” element-wisely performs the square-root operation if the input argument is a vector or a matrix. With  $\tilde{M} = \tilde{M}^T$ , the Hessian-vector product in CG to solve the diagonally-preconditioned linear system is

$$\hat{H}\hat{s} = \tilde{M}^{-1}(\tilde{H}(\tilde{M}^{-T}\hat{s})) = \tilde{M}^{-1}(\tilde{H}(\tilde{M}^{-1}\hat{s})). \quad (61)$$

To illustrate the computational details of PCG, we consider (27)’s Newton linear system as an example. First, we discuss the construction of the preconditioner. Consider the sub-problem (27) as an example. From

$$\tilde{\mathbf{x}}_i = \frac{1}{2}(\mathbf{x}_i \otimes \mathbf{q}_i) = \frac{1}{2}\text{vec}(\mathbf{q}_i \mathbf{x}_i^T),$$

we have

$$\begin{aligned}
& \sum_{i=1}^l D_{ii} (\tilde{\mathbf{x}}_i .* \tilde{\mathbf{x}}_i) \\
&= \frac{1}{4} \sum_{i=1}^l D_{ii} \text{vec} \left( (\mathbf{q}_i .* \mathbf{q}_i) (\mathbf{x}_i .* \mathbf{x}_i)^T \right) \\
&= \frac{1}{4} \sum_{i=1}^l \text{vec} \left( (\mathbf{q}_i .* \mathbf{q}_i) D_{ii} (\mathbf{x}_i .* \mathbf{x}_i)^T \right) \\
&= \frac{1}{4} \text{vec} ((Q .* Q) D (X .* X)).
\end{aligned} \tag{62}$$

Thus, the preconditioner of  $U$ 's sub-problem can be obtained via

$$\tilde{M} = \tilde{M}^T = \text{diag} \left( \sqrt{\lambda' \mathbf{1}_{nd \times 1} + \frac{1}{4} \text{vec} ((Q .* Q) D (X .* X))} \right)$$

or without vectorization

$$M = \sqrt{\lambda' \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} + \frac{1}{4} (Q .* Q) D (X .* X)} \in \mathcal{R}^{d \times n}. \tag{63}$$

We point out that the cost to obtain (63) is  $\mathcal{O}(d \times (\# \text{ nnz}))$ .

Next, we look into the computation of (61). Let  $\hat{S} \in \mathcal{R}^{d \times n}$  be the matrix form of the vector  $\hat{s}$  such that  $\hat{s} = \text{vec}(\hat{S})$ . It is possible to present (61) in terms of  $\hat{S}$  to avoid vectorizations. First, the matrix form of  $M^{-1} \hat{s}$  is  $\hat{S} ./ M$ , where “./” denotes the element-wise division. By substituting  $S = \hat{S} ./ M$  into (50) and then element-wisely dividing the resulting matrix by  $M$  again, (61) can be written as

$$\left( \lambda' (\hat{S} ./ M) + \frac{1}{2} Q \text{diag} (D \hat{\mathbf{z}}) X \right) ./ M, \tag{64}$$

where

$$\hat{\mathbf{z}} = \frac{1}{2} \left( Q^T .* \left( X (\hat{S} ./ M)^T \right) \right) \mathbf{1}_{d \times 1}.$$

For the computational complexity, besides  $\mathcal{O}(d \times (\# \text{ nnz}))$  from what we discussed in (50), the newly introduced divisions associated with  $M$  require  $\mathcal{O}(nd)$  operations. Unless the data matrix is very sparse, from  $n \leq \# \text{ nnz}$ , each PCG iteration does not cost significantly more than that without preconditioning. Note that we only compute  $M$  once in the beginning of PCG and then reuse it in each Hessian-vector product. See Algorithm 7 for a PCG procedure of using matrix variables and operations for solving the Newton system of  $U$ 's sub-problem.

#### 4.2. Sub-sampled Hessian Matrix

Some recent works have demonstrated that we can use sub-sampled Hessian to accelerate the Hessian-vector product in truncated Newton methods for empirical risk minimization [Byrd et al. 2011; Byrd et al. 2012; Wang et al. 2015]. From a statistical perspective, the training instances are drawn from an underlying distribution  $\Pr(y, \mathbf{x})$ .

---

**Algorithm 7** A preconditioned conjugate gradient method for solving (36) by operations on matrix variables. Linear systems in the sub-problem of  $U$  are considered.

---

- 1: Given  $0 < \eta < 1$  and  $G$ , the gradient matrix of the sub-problem (i.e., the matrix form of (47)). Let  $\hat{S} = \mathbf{0}_{d \times n}$ .
  - 2: Compute  $M$  via (63).
  - 3: Calculate  $R = -G ./ M$ ,  $\hat{D} = R$ , and  $\gamma^0 = \gamma = \|R\|_F^2$ .
  - 4: **while**  $\sqrt{\gamma} > \eta \sqrt{\gamma^0}$  **do**
  - 5:    $\hat{D}_h \leftarrow \hat{D} ./ M$
  - 6:    $\hat{z} \leftarrow (Q^T .* (X \hat{D}_h^T)) \mathbf{1}_{d \times 1}$
  - 7:    $\hat{D}_h \leftarrow (\lambda' \hat{D}_h + Q \text{diag}(D \hat{z}) X) ./ M$
  - 8:    $\alpha \leftarrow \gamma / \langle \hat{D}, \hat{D}_h \rangle$
  - 9:    $\hat{S} \leftarrow \hat{S} + \alpha \hat{D}$
  - 10:    $R \leftarrow R - \alpha \hat{D}_h$
  - 11:    $\gamma^{\text{new}} \leftarrow \|R\|_F^2$
  - 12:    $\beta \leftarrow \gamma^{\text{new}} / \gamma$
  - 13:    $\hat{D} \leftarrow R + \beta \hat{D}$
  - 14:    $\gamma \leftarrow \gamma^{\text{new}}$
  - 15: **end while**
  - 16: Output  $S = \hat{S} ./ M$  as the solution.
- 

Then, (34) can be interpreted as an empirical estimation of

$$\tilde{\lambda}I + lE [\xi''(\tilde{y}; y) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^T].$$

If  $L$  is a set of  $|L|$  instances randomly drawn from  $\{1, \dots, l\}$ , then the sub-sampled Hessian matrix,

$$\begin{aligned} \tilde{H}_{\text{sub}} &= \tilde{\lambda}I + \frac{l}{|L|} \sum_{i \in L} D_{ii} \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T, \\ &= \tilde{\lambda}I + \frac{l}{|L|} \tilde{X}_{L,:}^T D_{L,L} \tilde{X}_{L,:}, \end{aligned} \tag{65}$$

is an unbiased estimation of the original Hessian matrix. Note that the rows of  $\tilde{X}_{L,:} \in \mathcal{R}^{|L| \times \tilde{n}}$  are the sampled feature vectors and  $D_{L,L} \in \mathcal{R}^{|L| \times |L|}$  is the corresponding sub-matrix of  $D$ . After replacing  $\tilde{H}$  with  $\tilde{H}_{\text{sub}}$ , the linear system solved by CG becomes

$$\tilde{H}_{\text{sub}} \tilde{\mathbf{s}} = -\tilde{\mathbf{g}}. \tag{66}$$

Consider the sub-problem (27) as an example. From (50) and (54), if a subset of the rows in the data matrix  $X$  can be easily extracted, then the Hessian-vector product (in matrix form) for solving (66) can be conducted by

$$\lambda' S + \frac{l}{2|L|} Q_{:,L} \text{diag}(D_{L,L} \mathbf{z}_L) X_{L,:},$$

where  $Q_{:,L} \in \mathcal{R}^{d \times |L|}$  is a sub-matrix of  $Q$  containing  $q_i$ ,  $i \in L$  and

$$\mathbf{z}_L = \frac{1}{2} (Q_{:,L}^T .* (X_{L,:} S^T)) \mathbf{1}_{d \times 1}. \tag{67}$$

Therefore, a row-wise format should be used in storing  $X$ . In comparison with using all training instances, the cost of each CG iteration can be reduced to

$$\mathcal{O}\left(\frac{|L|d}{l} \times (\# \text{ nnz})\right),$$

when, for example, the sub-problem of  $U$  is considered. Although a small  $L$  can largely speed up the Hessian-vector product, the information loss caused by dropping instances can have a negative impact on the update direction found by solving (66). In the extreme case that  $L$  is an empty set, our truncated Newton method is reduced to a gradient descent method, which needs much more iterations. It may not be easy to decide the size of the subset  $L$ , but in Section 6.4 we will examine the running time under different choices.

## 5. PARALLELIZATION OF TRUNCATED NEWTON METHODS FOR SOLVING THE SUB-PROBLEMS

When solving each sub-problem, we should parallelize the main computational bottleneck: the calculations of the objective function, the gradient, and the Hessian-vector product. Because we have shown in Section 3.1 that each sub-problem is equivalent to a linear classification problem, our approach follows past developments of parallel Newton methods for regularized empirical risk minimization (e.g., [Zhuang et al. 2015] for distributed environments and [Lee et al. 2015] for multi-core environments). Here we focus on an implementation for multi-core machines by assuming that  $X$  can be instance-wisely accessed by multiple threads. To illustrate the details, we consider the sub-problem of  $U$  as an example. For easy description, we assume that all  $x_1, \dots, x_l$  are used without the sub-sampling technique in Section 4.2. If there are  $s$  threads available, we divide  $\{1, \dots, l\}$  into  $s$  disjoint subsets,  $L_1, \dots, L_s$ , for our task distribution. Recall the sub-matrix notations of  $X$ ,  $Q$ ,  $D$ , and  $z$  in Section 4.2. From (45) and (47), the parallelized forms of the output values and the gradient respectively are

$$\tilde{\mathbf{y}} = \begin{bmatrix} \frac{1}{2} (Q_{:,L_1}^T * (X_{L_1,:} U^T)) \mathbf{1}_{d \times 1} + X_{L_1,:} \mathbf{w} \\ \vdots \\ \frac{1}{2} (Q_{:,L_s}^T * (X_{L_s,:} U^T)) \mathbf{1}_{d \times 1} + X_{L_s,:} \mathbf{w} \end{bmatrix} \quad (68)$$

$$\nabla \tilde{f}(U) = \lambda' U + \frac{1}{2} \sum_{r=1}^s (Q_{:,L_r} \text{diag}(\mathbf{b}_{L_r}) X_{L_r,:}), \quad (69)$$

where tasks indexed by  $L_r$  is assigned to the  $r$ th thread. Form (50), parallel Hessian-vector product can be done by considering

$$\lambda' S + \frac{1}{4} \sum_{r=1}^s (Q_{:,L_r} \text{diag}(D_{L_r} \mathbf{z}_{L_r}) X_{L_r,:}). \quad (70)$$

For PCG, the diagonal preconditioner can be obtained via

$$M = \sqrt{\lambda' \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix} + \frac{1}{4} \sum_{r=1}^s (Q_{:,L_r} * Q_{:,L_r}) D_{L_r} (X_{L_r,:} * X_{L_r,:})} \in \mathcal{R}^{d \times n}. \quad (71)$$

Besides (68), for completing any of (69), (70), and (71), we must collect the results from the used threads by a series of  $d$ -by- $n$  matrix additions (called a  $dn$ -dimensional reduction) after the parallelized sections are finished. For example, if two threads are

Table II: Data statistics and parameters used in experiments. Density is the average number of non-zero features per instance.

Data set	$l$	$n$	nnz	$\lambda$	$\lambda'$	$\lambda''$
a9a	26,049	122	451,592	64	1	1
webspam	280,000	246	23,817,071	0.25	0.0625	0.0625
kddb	9,464,836	651,166	173,376,873	1	16	16
news20	16,009	1,355,191	7,303,887	0.0625	0.0625	0.0625
rcv1	677,399	47,236	39,627,321	0.0625	1	1
url	1,916,904	3,231,961	2,216,68,617	16	1	1
avazu-app	14,596,137	1,000,000	189,632,716	4	4	4
avazu-site	25,832,830	1,000,000	353,517,436	1	4	4
criteo	45,840,617	1,000,000	1,552,191,209	0.25	16	16

used together to compute (70), we can compute

$$B_1 = Q_{:,L_1} \text{diag}(D_{L_1} z_{L_1}) X_{L_1,:} \text{ and } B_2 = Q_{:,L_2} \text{diag}(D_{L_2} z_{L_2}) X_{L_2,:},$$

in parallel, and sum up the two matrices to obtain the second term of (70). Notice that from (51) and (70), the  $r$ th worker only maintains  $Q_{:,L_r}$  and  $z_{:,L_r}$  instead of the whole  $Q$  and  $z$ .

The parallelization of solving  $V$ 's sub-problem is very similar.

## 6. EXPERIMENTS

In some previous studies on linear regression and classification, coordinate descent methods beat Newton-type methods when the loss-related evaluations are not expensive [Chang et al. 2008; Ho and Lin 2012]. Since (29) is actually a linear regression or classification problem, it is not difficult to infer that ANT is not a good solver candidate to (4). Thus, we concentrate on (3) in the subsequent discussion and recommend using CD (or SG) otherwise.

To examine the effectiveness of ANT, we conduct a series of experiments on some real-world data sets. Section 6.1 gives the data statistics and discusses our experimental settings including the parameter selection. The environment for our experiments is described in Section 6.2. Sections 6.3-6.4 demonstrate the usefulness of the two techniques discussed in Section 4. The running time comparison on ANT against ADAGRAD and CD is presented in Section 6.5.

### 6.1. Data Sets and Parameters

We select several real-world data sets listed in Table II. All data sets can be downloaded at LIBSVM data sets page.<sup>2</sup> For webspam we consider the uni-gram version. For kddb, we choose the raw version of "bridge to algebra" track in KDD-Cup 2010 following [Juan et al. 2016]. To select proper parameters and calculate a test score, we need training, validation, and test sets for each problem. If a test set is available from the source, we directly use it.<sup>3</sup> Otherwise, we randomly select 20% instances from the whole data set for testing. Then, the training and validation sets are created by a 80-20 split of instances not used for testing. Note that we swap rcv1's training and test sets on LIBSVM page to make the training set bigger.

<sup>2</sup>LIBSVM data sets is available at [www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets](http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets).

<sup>3</sup>The labels in the test sets of avazu-app, avazu-site, and criteo are not meaningful, so they are treated as if they don't have test sets.

In problem (6), four parameters must be tuned. To avoid the expensive cost of searching the best setting, we fix  $d = 20$  and  $\lambda' = \lambda''$  throughout our experiments. For the regularization coefficients of each data set, we consider  $\{0.0625, 0.25, 1, 4, 16, 64\}$  as the search range of each regularization parameter and check all the combinations. Values leading to the lowest logistic loss on the validation set are used. Table II gives the parameters identified by the process. Furthermore, we set  $w_{\text{init}} = 0$  and randomly choose every element of  $U_{\text{init}}$  and  $V_{\text{init}}$  uniformly from  $[-1/\sqrt{d}, 1/\sqrt{d}]$ . For any comparison between different settings or algorithms, we ensure that they start with the same initial point. For the stopping condition of PCG, we set  $\eta = 0.3$  so Algorithm 7 terminates once

$$\|\hat{H}\hat{s} + \hat{g}\| \leq \eta\|\hat{g}\|.$$

Moreover, the two parameters of line search are  $\beta = 0.5$  and  $\nu = 0.01$ .

## 6.2. Environment and Implementation

The experiments in Sections 6.3-6.4 are conducted on a Linux machine with one Intel Core i7-6700 CPU 3.40GHz and 32 GB memory. For the experiments on multi-core implementations in Section 6.5, we migrate to another Linux machine with two Intel Xeon E5-2620 2.0GHz processors (12 physical cores in total) and 128 GB memory.

All the algorithms are implemented in C++ for a fair comparison. For the types of the variables, we use double-precision floating point for all real numbers and unsigned long integer for all indexes. We use `parallel-for` in OpenMP to implement all parallel computations discussed in Section 5. If the result generated by the threads is a scalar, the built-in scalar reduction can be used to automatically aggregate the results from different threads. For vector results, we manually merge them by sequential vector additions. This strategy has been confirmed to be effective in the parallelization of a truncated Newton method in [Lee et al. 2015]. For the parallelization of ADAGRAD, our implementation follows [Niu et al. 2011] to implement the for-loop in Algorithm 1 via `parallel-for` without the synchronization between them. The parallelization of CD is more difficult and beyond the scope of this work. Thus, we only implement a single-thread version.

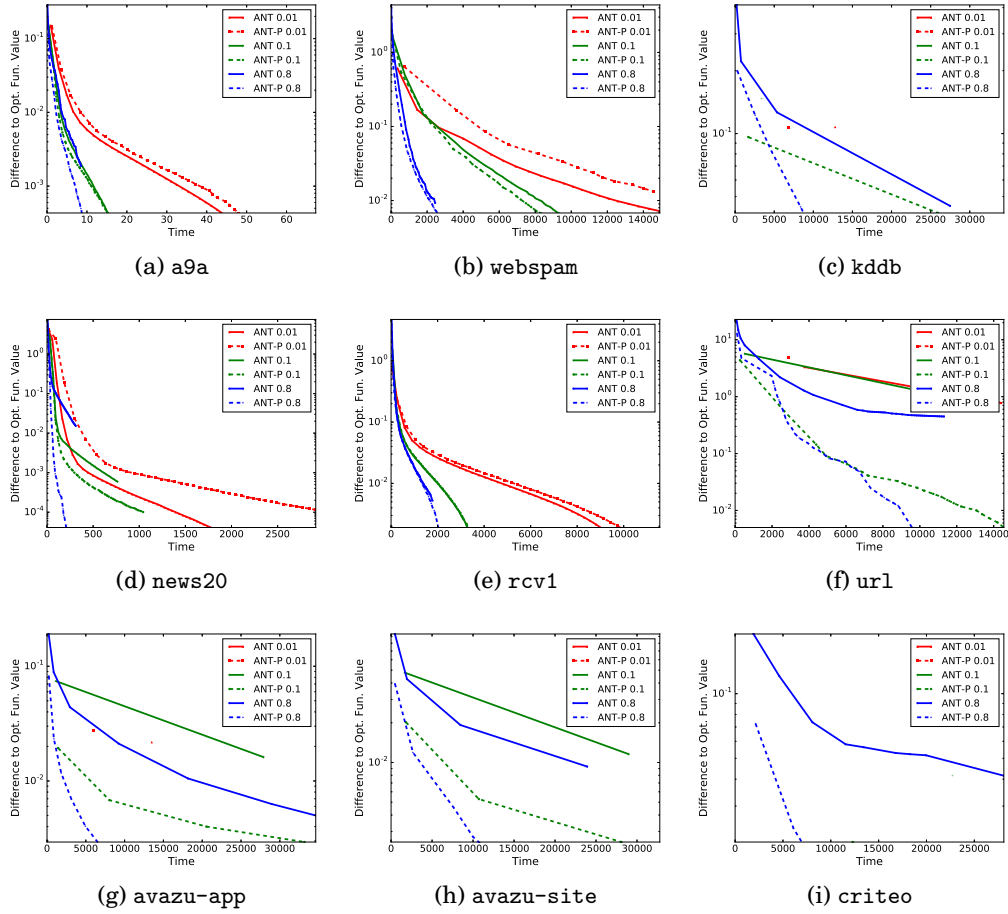
## 6.3. Stopping Tolerance of Sub-problems and the Effect of Preconditioning

To see the effect of different stopping tolerances in Algorithm 5, we consider  $\tilde{\epsilon} = 0.01, 0.1$ , and  $0.8$ , and examine the convergence speed of ANT. We also provide the results of ANT with PCG (denoted as ANT-P) to check the effectiveness of preconditioning. We present in Figure (1) the relation of running time to log-scaled distance between the current solution and a local optimal function value.

$$\frac{F(w, U, V) - F^*}{F^*}, \quad (72)$$

where  $F^*$  is the lowest function value reached among all settings.

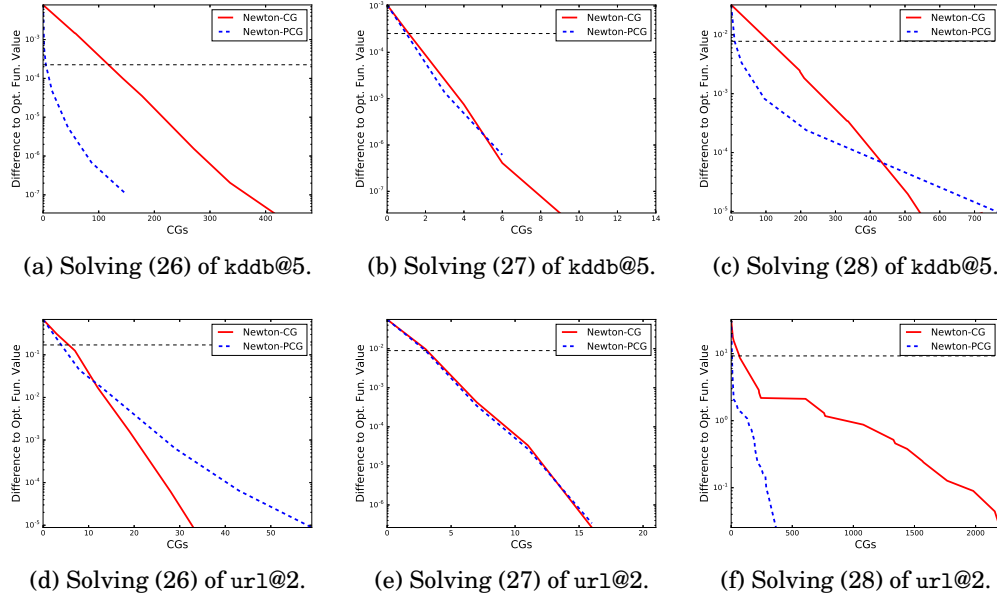
We begin with analyzing the results without preconditioning. From Figure 1, a larger inner stopping tolerance  $\tilde{\epsilon}$  leads to shorter overall training time. This result indicates that in the early stage of the alternating minimization procedure there is no need to waste time for accurately solving the sub-problems. For example, if neither  $V$  nor  $w$  is close to the optimum, in the sub-problem of  $U$ , getting an accurate solution is not very useful. Some earlier alternating minimization procedures have had similar observations and made the inner stopping condition from a loose one in the beginning to a tight one in the end; see, for example, [Lin 2007, Section 6]. Our relative stopping condition in (58) automatically achieves that; although  $\tilde{\epsilon}$  is fixed in the entire procedure, (58) is a strict condition in the final stage of the optimization process because of a small  $\|\nabla \tilde{f}(\tilde{w})\|$ .



**Fig. 1: Effects of the stopping condition and preconditioning for solving the sub-problems.** The number after ANT(P) indicates the used  $\tilde{\epsilon}$ . Time is in seconds. The  $y$ -axis is the log-scaled distance to a local optimal objective value; see (72). We terminate the training process if it does not produce any point within eight hours, so for larger data sets like criteo, some settings may have no line.

By comparing the ANT and ANT-P at different stopping conditions in Figure 1, we see that preconditioning is almost useless when the tolerance is tight (e.g.,  $\tilde{\epsilon} = 0.01$ ), but ANT-P converges faster than or at least comparable to ANT in the cases with a larger tolerance. To get more details, we investigate the convergence of the Newton method with/without preconditioning in solving sub-problems. We run ANT with  $\tilde{\epsilon} = 0.8$  and extract the three convex sub-problems (26)-(28) at a specific iteration. Then for each sub-problem we run Newton methods with/without preconditioning (denoted as Newton-PCG and Newton-CG, respectively), and present in Figure 2 the relation between the function-value reduction (of the sub-problem) and the number of CG iterations used. We have the following observations. First, in the beginning of solving each sub-problem, Newton-PCG performs as good as or even better than Newton-CG. Second, as the number of CG iterations increases, the performance gap between Newton-





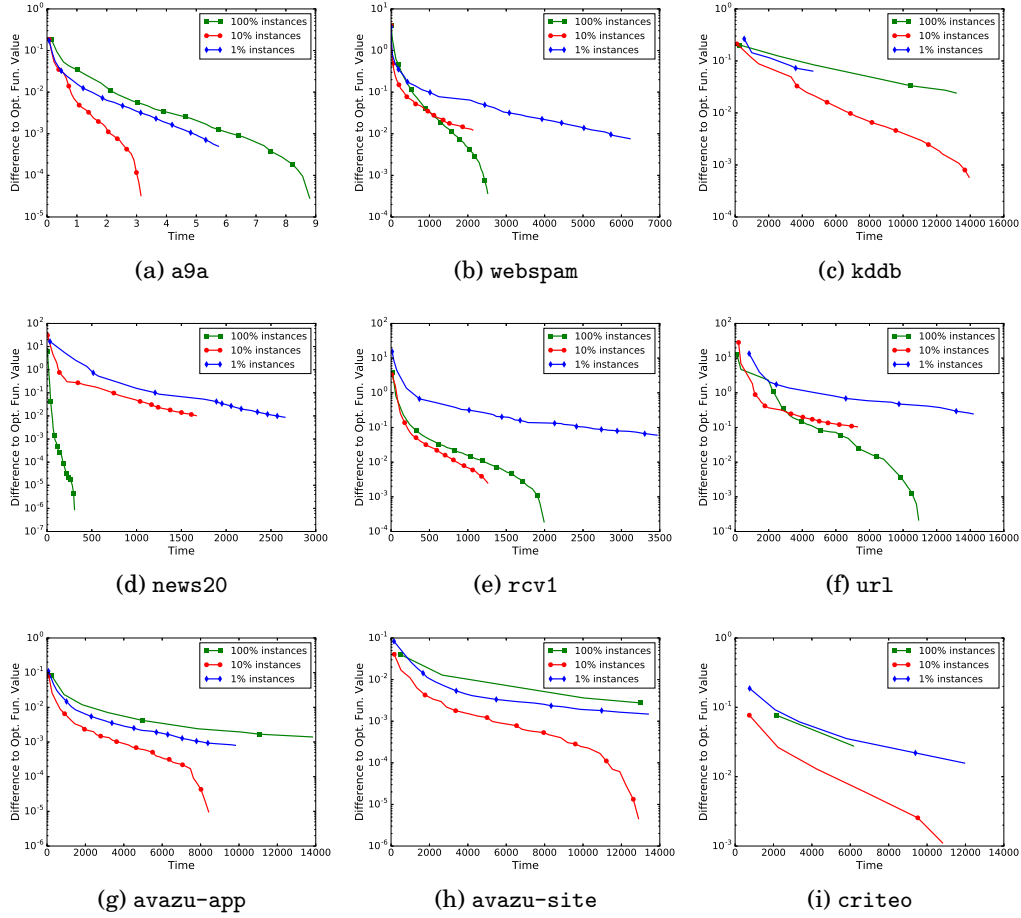
**Fig. 2: Analysis of Newton methods with/without preconditioning in solving the three sub-problems in an ANT iteration indicated by the integer after @.** The  $y$ -axis is the log-scaled distance to the global optimal function value of the sub-problem. The horizontal line indicates the first time that (58) with  $\epsilon = 0.8$  is satisfied during the optimization process. That is, in ANT (or ANT-P) the sub-problem is approximately solved until reaching the horizontal line.

CG and Newton-PCG gradually shrinks and preconditioning may even be harmful. Consequently, because we always terminate our truncated Newton method by a loose stopping of using  $\tilde{\epsilon} = 0.8$ , preconditioning is helpful or at least not harmful in ANT. A note that conducts more detailed analysis on the results in [Lin et al. 2008] has a similar conclusion.<sup>4</sup>

#### 6.4. Effects of the Use of Sub-sampled Hessian Matrix

In Section 4, we have learned that sub-sampled Hessian can save some operations in Hessian-vector products. However, the overall convergence may not be faster because of less accurate directions. To have a good understanding, we compare the results using the full Hessian matrix, a 10% sub-sampled Hessian matrix, and a 1% sub-sampled one by using the best setting in Section 6.3 (i.e., ANT-P with  $\tilde{\epsilon} = 0.8$ ). In Figure 3, for each sampling ratio, we show the change of objective values versus time. Comparing to using the full set, we observe that using 10% instances generally leads to similar or shorter training time. The difference is significant for data sets a9a, kddb, avazu-app, avazu-site, and criteo. These sets satisfy  $l \gg n$ , so some instances may carry redundant information and can be dropped. However, when only 1% instances are used, ANT often spends more time to achieve a specified objective value than that of using the full set.

<sup>4</sup>[www.csie.ntu.edu.tw/~cjlin/papers/logistic/pcgnote.pdf](http://www.csie.ntu.edu.tw/~cjlin/papers/logistic/pcgnote.pdf)



**Fig. 3: Effects of ANT-P with different rates of using sub-sampled Hessian.** Time is in seconds. The  $y$ -axis is the log-scaled distance to a local optimal objective value; see (72).

We conclude that by selecting a suitable amount of instances, the approach of using sub-sampled Hessian is useful for ANT. However, it is important to avoid choosing a too small subset.

### 6.5. A Comparison on ANT and Other State of the Art Methods

We compare the running time of ANT-P against ADAGRAD and CD described in Section 2 under both single-thread and multi-thread settings. The parameters of ANT-P is determined by our discussion in Sections 6.3 and 6.4; we choose  $\tilde{\epsilon} = 0.8$  as the stopping condition of solving sub-problems and use 10% of data for constructing the sub-Hessian matrix. For the parameter  $\eta_0$  in the stochastic gradient method detailed in Algorithm 1, we use 0.1. For CD's number of inner iterations, we set  $\bar{k}_{inner}$  for 1. Following the format of Figure 1, we draw these methods' running performance at different numbers of threads in Figure 4. Note that the parallelization of CD is not easy, so we only implement a single-thread version. From Figure 4, we see that CD is largely outperformed

Table III: Portion of time spent on  $\exp / \log$  operations.

Solver \ Data set	a9a	news20
ADAGRAD	0.64%	0.06%
CD	81.50%	65.05%
ANT	3.52%	0.07%

Table IV: Test logistic losses of LR and the two FM formulations (2) and (6).

	Model	LR	FM: (2)	FM: (6)		
	Solver	Algorithm 5	ADAGRAD	ADAGRAD	CD	ANT-P
Data set	a9a	0.3238	0.3200	0.3200	0.3206	0.3204
	webspam	0.2085	0.1151	0.0904	NA	0.0583
	news20	0.1194	0.1019	0.0879	0.0824	0.0844
	kddb	0.2852	0.2652	0.2647	NA	0.2673
	url	0.0315	0.0144	0.0155	NA	0.0157
	rcv1	0.0740	0.0623	0.0517	0.0506	0.0505
	avazu-app	0.3392	0.3407	0.3300	NA	0.3344
	avazu-site	0.4460	0.4377	0.4354	NA	0.4381
	criteo	0.4616	0.4474	0.4489	NA	0.4513

by ADAGRAD and ANT-P. From Sections 2.1, 2.2, and 3.3, a possible reason is that CD costs much more expensive loss-related operations (i.e.,  $\exp / \log$  evaluations) than its counterparts. To support our argument, table III gives the time portion of executing  $\exp / \log$  functions in the first ANT-P iteration, the first epoch of ADAGRAD, and the first outer iteration of CD on a9a and news20. This observation is consistent with a study on LR [Yuan et al. 2010]. Furthermore, we observe that the parallelization of ANT-P and ADAGRAD is effective. Their training time decreases as the number of threads increases. On the other hand, in either single- and multi-thread settings, ANT-P always produces significantly smaller objective values comparing with ADAGRAD and CD. Two possible reasons are given. One is that ADAGRAD may eventually reach the same values but its convergence is very slow. The other one is that the FM problem is not convex, so ADAGRAD and ANT may reach different solutions.

## 6.6. Comparison on Test Performance with Other Formulations

From a practical perspective, it is important to examine if the modified FM problem (6) can perform as well as the original FM problem (2). To this end, for each data set, we compare their test accuracies and logistic losses in Tables IV-V. To see how FM improves over a linear model, we include results of LR. For a fair comparison, we use  $d = 40$  for (1) and  $d = 20$  for (5) so that they have the same number of model parameters. For the solvers, we consider the truncated Newton method in Algorithm 5 for LR, the stochastic gradient method in Section 2.1 for (2) and (6), and CD and ANT-P for (6). As shown in Tables IV-V, the performance gap between problems (2) and (6) is not significant, so the modified FM formulation in (6) seems to give similar learning capability. Regarding the comparison between FM and LR, FM gives a much higher test score than LR on webspam. This data set has less features than its instances, so feature conjunction seems to be useful. For other data sets, the performance gap is smaller. For a9a, it is known that even highly non-linear Gaussian-kernel SVM gives

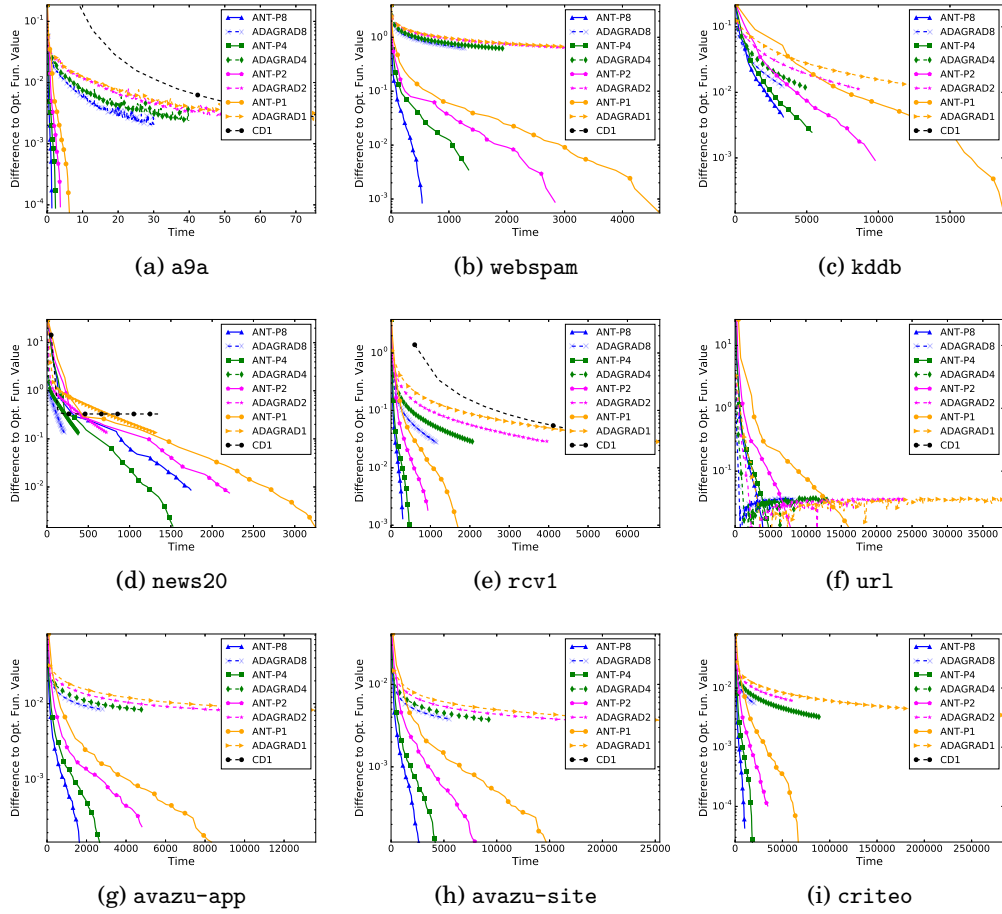


Fig. 4: A comparison between ANT-P, CD, and ADAGRAD. The number in the legend indicates the number of threads. Time is in seconds. The  $y$ -axis is the log-scaled distance to a local optimal objective value; see (72).

only similar accuracy to a linear classifier [Huang and Lin 2016]. For the seven sparse data sets, three of them are document data, and generally linear classifiers are powerful enough [Yuan et al. 2012b]. However, overall we still see that FM slightly boosts the test accuracy and logistic loss.

One may think that the test scores in Table IV-V do not reflect the gap between ANT-P and ADAGRAD in Figure 4. Again, two possible reasons are given. One is that ADAGRAD may eventually reach the same values but its convergence is very slow. The other one is that the FM problem is not convex, so ADAGRAD and ANT may reach different solutions.

## 7. CONCLUSIONS

In this paper, we propose an effective alternating Newton method for solving a multi-block convex reformulation of factorization machines. A sequence of convex sub-problems are solved, where we show that each sub-problem can be written in a form

Table V: Test accuracies of LR and the two FM formulations (2) and (6).

	Model	LR	FM: (2)	FM: (6)		
	Solver	Algorithm 5	ADAGRAD	ADAGRAD	CD	ANT-P
Data set	a9a	85.03%	85.26%	85.24%	85.21%	85.18%
	webspam	92.63%	95.93%	97.00%	NA	98.34%
	news20	96.43%	95.98%	96.38%	96.78%	96.61%
	kdddb	88.87%	89.43%	89.62%	NA	89.52%
	url	98.87%	99.53%	99.51%	NA	99.50%
	rcv1	97.59%	98.05%	98.39%	98.38%	98.39%
	avazu-app	87.09%	87.11%	87.13%	NA	87.06%
	avazu-site	80.44%	80.61%	80.67%	NA	80.58%
	criteo	78.37%	79.13%	79.04%	NA	78.91%

equivalent to regularized logistic regression. Because variables in a sub-problem may be a matrix rather than a vector, we carefully design a truncated Newton method that employs fast matrix operations. Further, some advanced techniques such as preconditioned conjugate gradient methods and sub-sampled Hessian technique are incorporated into our framework. Through experiments in this paper, we establish the superiority over existing stochastic gradient algorithms and coordinate descent methods. We also show the effectiveness of the parallelization of our algorithm.

In summary, we have successfully developed a useful algorithm and implementation for training factorization machines. A MATLAB package of the proposed method and programs used for experiments are available at

<http://www.csie.ntu.edu.tw/~cjlin/papers/fm>.

With the built-in matrix operations, our MATLAB package is as efficient as our C++ implementation but the MATLAB training and prediction modules only need around 150 lines of code. It means that our algorithm can be easily ported into any existing service as long as a good matrix library is available.

## APPENDIX

Here we establish the convergence for cyclic block coordinate descent (CBCD) when solving (6) so that the convergence of ANT and CD automatically holds. Consider an unconstrained optimization problem

$$\min_{\mathbf{w} \in \mathcal{R}^n} f(\mathbf{w}). \quad (73)$$

First, we assume that  $\mathbf{w}$  are divided into  $B$  blocks. To get an update direction  $\mathbf{s}_b$  for the  $b$ th block, we approximately solve

$$\min_{\mathbf{s}_b \in \mathcal{R}^n} \nabla_b f(\mathbf{w})^T \mathbf{s}_b + \frac{1}{2} \mathbf{s}_b^T \nabla_b^2 f(\mathbf{w}) \mathbf{s}_b, \quad (74)$$

where the elements in  $\mathbf{s}_b$ ,  $\nabla_b f(\mathbf{w}_b)$ , and  $\nabla_b^2 f(\mathbf{w})$  are set to zero if they are not associated with the  $b$ th block. Then, back-tracking line search is used to find a step  $\theta$  and the solution is changed to

$$\mathbf{w} + \theta \mathbf{s}_b.$$

Let  $\mathbf{w}_b^k$  denote the initial value of  $\mathbf{w}$  in the  $k$ th update of the  $b$ th block. Algorithm 8 shows CBCD in detail.

**Algorithm 8** Cyclic block coordinate descent method.

---

```

1: Given an initial solution  $w_1^1$ .
2: for  $k \leftarrow \{1, \dots\}$  do
3:   for  $b \leftarrow \{1, \dots, B\}$  do
4:     Find a update direction  $s_b^k$  by solving (74).
5:     Conduct line search to find  $\theta_b^k$ .
6:      $w_{b+1}^k \leftarrow w_b^k + \theta_b^k s_b^k$ 
7:   end for
8: end for

```

---

ASSUMPTION .1. *The objective function  $f(w)$  is smooth and bounded below. Moreover, the magnitude of  $\nabla f(w)$  is upper bounded.*

Definition .2. For updating the block  $b$ , a vector  $s_b$  is called a descent direction if

$$-\nabla_b f(w)^T s_b > 0 \text{ when } \|\nabla_b f(w)\| > 0.$$

THEOREM .3. *If  $s$  is a descent direction, back-tracking line search must end up in a finite number of iterations. That is,  $\exists \theta > 0$ , s.t.*

$$f(w + \theta s) \leq f(w) + \eta \theta \nabla f(w)^T s, \quad (75)$$

where  $0 < \eta < 1$ .

PROOF. This theorem can be proved via contradiction. Assume that the stopping condition of line search is never satisfied. Thus, for  $\theta \in \{1, \beta, \beta^2, \dots\}$ , we have

$$\frac{f(w + \theta s) - f(w)}{\theta} > \eta \nabla f(w)^T s.$$

It leads to a contradiction because

$$\lim_{\theta \rightarrow 0} \frac{f(w + \theta s) - f(w)}{\theta} = \nabla f(w)^T s > \eta \nabla f(w)^T s$$

is wrong. Q.E.D.  $\square$

THEOREM .4. *If the (sub-)Hessian matrix of each sub-problem satisfies*

$$\sigma_1 I \preceq \nabla_b^2 f(w) \preceq \sigma_2 I, \quad (76)$$

*except in the first CG iteration, every CG iterate meets*

$$-\nabla_b f(w)^T s \geq c_1 \|s\|^2 \text{ and } \|s\| \geq c_2 \|\nabla_b f(w)\|. \quad (77)$$

PROOF.  $\square$

THEOREM .5. *Every clustering point of  $\{w_b^k\}$  is a stationary point.*

PROOF. The use of line search in CBCD ensures that a function decrease in each step, so  $\{f(w_b^k)\}$  is a decreasing sequence or equivalently

$$\dots \geq f(w_1^k) \geq \dots \geq f(w_B^k) \geq f(w_1^{k+1}) \geq \dots \geq f(w_B^{k+1}) \geq \dots \quad (78)$$

Since the objective function is lower-bounded,  $\{f(w_b^k)\}$  must converge to a finite value  $f^*$ . Recalling  $w_{b+1}^k = w_b^k + \theta_b^k s_b^k$  and line search's stopping condition, we get

$$f(w_b^k) - f(w_{b+1}^k) \geq -\eta \theta_b^k \nabla_b f(w_b^k)^T s_b^k \geq 0. \quad (79)$$

Because  $\{f(w_b^k)\}$  is a convergent sequence, the left-hand side in (79) converges to 0 and therefore by Squeeze Theorem, we have

$$\{\theta_b^k \nabla_b f(w_b^k)^T s_b^k\} \rightarrow 0. \quad (80)$$

Assume  $\bar{w}$  is a non-stationary cluster point of  $\{w_b^k\}$ . There must be a convergent subsequence  $\{w_c^k\}_{\mathcal{K}} \rightarrow \bar{w}$  with infinitely many updates of a block  $c \in \{1, \dots, B\}$ . Next, we prove that  $\|\nabla_c f(\bar{w})\| = 0$  by contradiction. Suppose that  $\|\nabla_c f(\bar{w})\| > 0$ . By the smoothness of  $f(w)$ , we can pick up  $k_0$  large enough so that

$$\|\nabla_c f(\bar{w}_c^k)\| > \epsilon \text{ if } k > k_0. \quad (81)$$

With (80), (81) implies  $\{\theta_c^k\}_{\mathcal{K}_0} \rightarrow 0$ , so there must be an iteration index  $k_1$  large enough to guarantee that the Armijo condition does not hold in the beginning. Let  $\mathcal{K}_1 = \{k \mid k \in \mathcal{K}_0, k > k_1\}$ . For updating  $w_c^k, k \in \mathcal{K}_1$ , the last violation of Armijo rule indicates

$$f(w_c^k) - f(w_c^k + \hat{\theta}_c^k s_c^k) < -\eta \hat{\theta}_c^k \nabla_c f(w_c^k)^T s_c^k, \quad (82)$$

where  $\hat{\theta}_c^k = \frac{\theta_c^k}{\beta}$ . By applying Mean Value Theorem on the left-hand side of (82), we can rewrite (82) into

$$-\nabla_c f(w_c^k + \tilde{\theta}_c^k s_c^k)^T s_c^k \hat{\theta}_c^k < -\eta \hat{\theta}_c^k \nabla_c f(w_c^k)^T s_c^k. \quad (83)$$

where  $\tilde{\alpha}^{k,c} \in [0, \bar{\alpha}^{k,c}]$ . Because bounded  $\nabla f(x)$  means bounded  $s_c^k$ , we can find a subsequence of  $\mathcal{K}_1, \tilde{\mathcal{K}}_1$ , such that  $\{s_c^k\}_{\tilde{\mathcal{K}}_1}$  converges to a fixed point  $\bar{s}_c$ . Considering  $k \in \mathcal{K}_2^c$  and taking  $k \rightarrow \infty$  for (83), we obtain

$$(1 - \eta) \nabla_c f(\bar{w})^T \bar{s}_c > 0. \quad (84)$$

However, (84) contradicts descent property of the search direction. Thus,  $\nabla_c f(\bar{w}) = 0$ .

Because  $\nabla_c f(\bar{w}) = 0, w_{c+1}^k \rightarrow w_c^k$  for  $k \in \mathcal{K}$ . That is,  $\{w_{c+1}^k\}_{\mathcal{K}}$  is also a convergent sequence with limit at  $\bar{w}$ . By an analogy of proving  $\nabla_c f(\bar{w}) = 0$ , we can also show  $\nabla_{c+1} f(\bar{w}) = 0$ . By looping the same strategy for all other remaining blocks, the proof is completed because the gradient at each block is zero.  $\square$

## REFERENCES

- Dimitri P. Bertsekas. 1999. *Nonlinear Programming* (second ed.). Athena Scientific, Belmont, MA 02178-9998.
- Mathieu Blondel, Masakazu Ishihata, Akinori Fujino, and Naonori Ueda. 2016. Polynomial networks and factorization machines: new insights and efficient training algorithms. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*.
- Richard H. Byrd, Gillian M. Chin, Will Neveitt, and Jorge Nocedal. 2011. On the use of stochastic Hessian information in optimization methods for machine learning. *SIAM Journal on Optimization* 21, 3 (2011), 977–995.
- Richard H. Byrd, Gillian M. Chin, Jorge Nocedal, and Yuchen Wu. 2012. Sample size selection in optimization methods for machine learning. *Mathematical Programming* 134, 1 (2012), 127–155.
- Kai-Wei Chang, Cho-Jui Hsieh, and Chih-Jen Lin. 2008. Coordinate descent method for large-scale L2-loss linear SVM. *Journal of Machine Learning Research* 9 (2008), 1369–1398. <http://www.csie.ntu.edu.tw/~cjlin/papers/cdl2.pdf>
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12 (2011), 2121–2159.
- Stanley C. Eisenstat and Homer F. Walker. 1994. Globally convergent inexact newton methods. *SIAM Journal on Optimization* 4, 2 (1994), 393–422.
- Gene H. Golub and Charles. F. Van Loan. 2012. *Matrix Computations* (fourth ed.). The Johns Hopkins University Press.
- Chia-Hua Ho and Chih-Jen Lin. 2012. Large-scale linear support vector regression. *Journal of Machine Learning Research* 13 (2012), 3323–3348. <http://www.csie.ntu.edu.tw/~cjlin/papers/linear-svr.pdf>

- Fang-Lan Huang, Cho-Jui Hsieh, Kai-Wei Chang, and Chih-Jen Lin. 2010. Iterative scaling and coordinate descent methods for maximum entropy. *Journal of Machine Learning Research* 11 (2010), 815–848. <http://www.csie.ntu.edu.tw/~cjlin/papers/maxent.journal.pdf>
- Hsin-Yuan Huang and Chih-Jen Lin. 2016. Linear and kernel classification: When to use which?. In *Proceedings of SIAM International Conference on Data Mining (SDM)*. <http://www.csie.ntu.edu.tw/~cjlin/papers/kernel-check/kcheck.pdf>
- Yuchin Juan, Yong Zhuang, Wei-Sheng Chin, and Chih-Jen Lin. 2016. Field-aware factorization machines for CTR prediction. In *Proceedings of the ACM Recommender Systems Conference (RecSys)*. <http://www.csie.ntu.edu.tw/~cjlin/papers/ffm.pdf>
- S. Sathya Keerthi and Dennis DeCoste. 2005. A modified finite Newton method for fast solution of large scale linear SVMs. *Journal of Machine Learning Research* 6 (2005), 341–361.
- Paul Komarek and Andrew W. Moore. 2005. *Making Logistic Regression A Core Data Mining Tool: A Practical Investigation of Accuracy, Speed, and Simplicity*. Technical Report TR-05-27. Robotics Institute, Carnegie Mellon University.
- Mu-Chu Lee, Wei-Lin Chiang, and Chih-Jen Lin. 2015. Fast matrix-vector multiplications for large-scale logistic regression on shared-memory systems. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*. <http://www.csie.ntu.edu.tw/~cjlin/papers/multicore.liblinear.icdm.pdf>
- Chih-Jen Lin. 2007. Projected gradient methods for non-negative matrix factorization. *Neural Computation* 19 (2007), 2756–2779. <http://www.csie.ntu.edu.tw/~cjlin/papers/pgadnmf.pdf>
- Chih-Jen Lin, Ruby C. Weng, and S. Sathya Keerthi. 2008. Trust region Newton method for large-scale logistic regression. *Journal of Machine Learning Research* 9 (2008), 627–650. <http://www.csie.ntu.edu.tw/~cjlin/papers/logistic.pdf>
- James Martens. 2010. Deep learning via Hessian-free optimization. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*.
- Stephen G. Nash. 2000. A survey of truncated-Newton methods. *Journal of Computational and Applied Mathematics* 124, 1–2 (2000), 45–59.
- Stephen G. Nash and Ariele Sofer. 1996. *Linear and Nonlinear Programming*. McGraw-Hill.
- Feng Niu, Benjamin Recht, Christopher Ré, and Stephen J. Wright. 2011. HOGWILD!: a lock-free approach to parallelizing stochastic gradient descent. In *Advances in Neural Information Processing Systems 24*, J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F.C.N. Pereira, and K.Q. Weinberger (Eds.). 693–701.
- Steffen Rendle. 2010. Factorization machines. In *Proceedings of IEEE International Conference on Data Mining (ICDM)*. 995–1000.
- Steffen Rendle. 2012. Factorization machines with libFM. *ACM Transactions on Intelligent Systems and Technology (TIST)* 3, 3 (2012), 57.
- Nicol N Schraudolph. 2002. Fast curvature matrix-vector products for second-order gradient descent. *Neural Computation* 14, 7 (2002), 1723–1738.
- Anh-Phuong Ta. 2015. Factorization machines with follow-the-regularized-leader for CTR prediction in display advertising. In *Proceedings of the IEEE International Conference on Big Data*.
- Chien-Chih Wang, Chun-Heng Huang, and Chih-Jen Lin. 2015. Subsampled Hessian Newton methods for supervised learning. *Neural Computation* 27 (2015), 1766–1795. <http://www.csie.ntu.edu.tw/~cjlin/papers/sub.hessian/sample.hessian.pdf>
- Guo-Xun Yuan, Kai-Wei Chang, Cho-Jui Hsieh, and Chih-Jen Lin. 2010. A comparison of optimization methods and software for large-scale l1-regularized linear classification. *Journal of Machine Learning Research* 11 (2010), 3183–3234. <http://www.csie.ntu.edu.tw/~cjlin/papers/l1.pdf>
- Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin. 2012a. An improved GLMNET for l1-regularized logistic regression. *Journal of Machine Learning Research* 13 (2012), 1999–2030. <http://www.csie.ntu.edu.tw/~cjlin/papers/l1.glmnet/long-glmnet.pdf>
- Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin. 2012b. Recent advances of large-scale linear classification. *Proc. IEEE* 100, 9 (2012), 2584–2603. <http://www.csie.ntu.edu.tw/~cjlin/papers/survey-linear.pdf>
- Yunhong Zhou, Dennis Wilkinson, Robert Schreiber, and Rong Pan. 2008. Large-scale parallel collaborative filtering for the Netflix prize. In *Proceedings of the Fourth International Conference on Algorithmic Aspects in Information and Management*. 337–348.
- Yong Zhuang, Wei-Sheng Chin, Yu-Chin Juan, and Chih-Jen Lin. 2015. Distributed Newton method for regularized logistic regression. In *Proceedings of the Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD)*.

Received ; revised ; accepted