



國立臺灣大學電機資訊學院資訊工程學研究所

碩士論文

Department of Graduate Institute of Computer Science and Information  
Engineering

College of Electrical Engineering and Computer Science

National Taiwan University

Master Dissertation

Newton Methods for Factorization Machines

黃郁庭

Yu-Ting Huang

指導教授：林智仁 博士

Advisor: Chih-Jen Lin, Ph.D.

中華民國 108 年 00 月

XX, 2019

## 中文摘要



...

關鍵詞: ...



## ABSTRACT

...

KEYWORDS:



## TABLE OF CONTENTS

中文摘要 . . . . .	i
ABSTRACT . . . . .	ii
LIST OF FIGURES . . . . .	iv
LIST OF TABLES . . . . .	v
CHAPTER	
I. Introduction . . . . .	1
II. Newton Method for Unconstrained Minimization . . . . .	3
III. Newton Methods for Factorization Machines . . . . .	5
3.1 Gradient Calculation . . . . .	6
3.2 Hessian and Gauss-Newton Matrices . . . . .	8
3.3 Implementation of Line Search . . . . .	9
3.4 Overall Procedure . . . . .	10
3.5 Implementation for Matrix Factorization . . . . .	12
IV. A Review of Alternating Newton Methods . . . . .	17
V. Experiments . . . . .	21
5.1 Experimental Settings . . . . .	21



## LIST OF FIGURES

### Figure

5.1	A comparison on the va-loss of different methods. ANT is worse than Gauss among smaller $\lambda$ . The $x$ -axis is the iteration, while the $y$ -axis is the va-loss. . . . .	22
5.2	A comparison on the va-loss of different methods. The $x$ -axis is the iteration, while the $y$ -axis is the va-loss. . . . .	23
5.3	ml-1m. A comparison on the va-loss of different $\lambda$ . The $x$ -axis is the iteration, while the $y$ -axis is the va-loss. . . . .	24
5.4	ml-1m. A comparison on the best va-loss of different methods. The $x$ -axis is $\lambda$ , while the $y$ -axis is the va-loss. . . . .	25



## LIST OF TABLES

### Table

5.1	Statistics and parameters for each data set. . . . .	21
-----	--	----



## CHAPTER I

### Introduction

Given a rating matrix  $R$ , matrix factorization (MF) is a technique to find two dense matrices  $U \in \mathcal{R}^{d \times M}$  and  $V \in \mathcal{R}^{d \times N}$  such that  $r_{m,n} \simeq \mathbf{u}_m^T \mathbf{v}_n$ , where  $\mathbf{u}_m \in \mathcal{R}^d$  and  $\mathbf{v}_n \in \mathcal{R}^d$  are respectively the  $m$ th column of  $U$  and the  $n$ th column of  $V$ ,  $d$  is the pre-specified number of latent features, and the entry  $r_{m,n}$  denotes the feedback of the  $m$ th user on the  $n$ th item. This task is achieved by solving the following non-convex problem

$$\min_{U,V} \frac{1}{2} \sum_{(m,n) \in R} (r_{m,n} - \mathbf{u}_m^T \mathbf{v}_n)^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) \quad (1.1)$$

where  $(m, n) \in R$  indicates that rating  $r_{m,n}$  is available,  $\lambda$  is regularization coefficients for avoiding over-fitting, and  $\|\cdot\|_F$  is the Frobenius norm. In fact, problem (1.1) is a special case of Factorization Machines (FM) (Rendle, 2010). The optimization problem solved by FM is as follows.

$$\min_{U,V} f(U, V) \quad (1.2)$$

where

$$f(U, V) = \frac{1}{2} \sum_{i=1}^l \left( y_i - (U \mathbf{w}_i)^T (V \mathbf{h}_i) \right)^2 + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) \quad (1.3)$$

and  $l$  is the total number of training instances. Problem (1.3) can be considered as a regression problem, where  $y_i$  is the label, and  $\mathbf{h}_i \in \mathcal{R}^N$  and  $\mathbf{w}_i \in \mathcal{R}^M$  are the feature vectors. If we let  $l$  be the number of non-zero entries in  $R$ , the value  $y_i = r_{m,n}$  be the corresponding

rating from  $R$ ,

$$\mathbf{w}_i = [\underbrace{0, \dots, 0}_{m-1}, 1, 0, \dots, 0]^T, \text{ and} \quad (1.4)$$

$$\mathbf{h}_i = [\underbrace{0, \dots, 0}_{n-1}, 1, 0, \dots, 0]^T, \quad (1.5)$$



then problem (1.2) is reduced to (1.1). Note that (1.2) is a variant of FM proposed in Blondel et al., (2016). One of its advantages is that the function in (1.3) is multi-block convex. That is, if  $U$  (or  $V$ ) is fixed, then  $f(U, V)$  is a convex function of  $V$  (or  $U$ ). In this work, we focus on solving the optimization problem (1.2).

Many optimization methods have been proposed to minimize (1.2). We are particularly interested in Newton methods. Chin et al., (2018) have proposed an alternating Newton method and showed that it achieves state-of-the-art efficiency. The basic idea is to alternately switch between solving two sub-problems, each of which minimizes over one block of variables but fixes the other. Namely the first sub-problem minimizes over  $U$  while keeping  $V$  fixed, and the other minimizes over  $V$  while keeping  $U$  fixed.

Instead of alternately solving two sub-problems, naturally we ask why not minimizing  $f(U, V)$  directly by the Newton method. To the best of our knowledge, no study has considered such a setting for factorization machines. In this work, we begin with developing a Newton method to solve problem (1.2). We then conduct detailed comparisons with the alternating Newton method.





## CHAPTER II

### Newton Method for Unconstrained Minimization

To derive the Newton method, we should re-write the function to have a vector of variables. To this end, the objective function in (1.3) is changed to

$$f(\text{vec}(U), \text{vec}(V)),$$

where  $\text{vec}(\cdot)$  stacks all columns of a matrix to a vector. For the standard Newton methods, at the  $k$ th iteration, we find a direction  $\mathbf{s}_k$  minimizing the following second-order approximation of the function value:

$$\min_{\mathbf{s}} \quad \frac{1}{2} \mathbf{s}^T H_k \mathbf{s} + \mathbf{g}_k^T \mathbf{s} \quad (2.1)$$

where

$$H_k = \nabla^2 f(\text{vec}(U_k), \text{vec}(V_k)) \text{ and } \mathbf{g}_k = \nabla f(\text{vec}(U_k), \text{vec}(V_k))$$

are the Hessian matrix and the gradient of  $f(U_k, V_k)$ , respectively. If  $H_k$  is positive definite, then (2.1) is equivalent to solving the following linear system:

$$H_k \mathbf{s} = -\mathbf{g}_k \quad (2.2)$$

For non-convex objective functions, the Hessian matrix may not be positive definite. Then the solution of (2.2) may not lead to a descent direction. In Section IV we will show that this issue occurs for our optimization problem (1.2). To obtain a descent direction, we can consider a positive-definite approximation of the Hessian matrix.

Now assume that  $G_k$  is a positive definite approximation of the Hessian matrix. Instead of solving (2.2), we solve the following linear system to find a direction  $\mathbf{s}_k$ .

$$G_k \mathbf{s} = -\mathbf{g}_k \quad (2.3)$$

In practice, (2.3) may not need to be exactly solved. Therefore, truncated Newton methods have been introduced to approximately solve (2.3). Iterative methods such as conjugate gradient method are often used for approximately solving (2.3), so at each iteration an inner iterative process is involved. In particular, if CG is used, at each inner iteration a matrix-vector product between  $G_k$  and a vector must be conducted.

After a direction is found in our Newton method, we must decide the step size taken along that direction. A common setting is the backtracking line search. Namely, we find the largest step size  $\theta \in \{1, \beta, \beta^2, \dots\}$  satisfying the following sufficient decrease condition. Let

$$\mathbf{s} = \text{vec}(S) \quad (2.4)$$

and

$$S = \begin{bmatrix} S_u & S_v \end{bmatrix}. \quad (2.5)$$

$$f(U + \theta S_u, V + \theta S_v) - f(U, V) \leq \theta \nu \mathbf{g}^T \text{vec}(S), \quad (2.6)$$

where  $\nu \in (0, 1)$  and  $\beta \in (0, 1)$  are pre-specified constants.



## CHAPTER III

### Newton Methods for Factorization Machines

To compute the gradient and Hessian of (1.2), we begin with the following properties of the Kronecker product among matrices  $A$ ,  $B$  and  $X$ :

$$(B^T \otimes A) \text{vec}(X) = \text{vec}(AXB) \quad (3.1)$$

$$(A \otimes B)^T = A^T \otimes B^T. \quad (3.2)$$

Moreover, the predicted value can be represented as

$$(U\mathbf{w}_i)^T(V\mathbf{h}_i) = (\mathbf{h}_i^T \otimes (U\mathbf{w}_i)^T) \text{vec}(V) \quad (3.3)$$

$$= (\mathbf{h}_i \otimes (U\mathbf{w}_i))^T \text{vec}(V), \quad (3.4)$$

where (3.3) and (3.4) are from (3.1) and (3.2), respectively. The predicted value also can be represented as

$$\begin{aligned} (U\mathbf{w}_i)^T(V\mathbf{h}_i) &= (V\mathbf{h}_i)^T(U\mathbf{w}_i) \\ &= (\mathbf{w}_i^T \otimes (V\mathbf{h}_i)^T) \text{vec}(U) \\ &= (\mathbf{w}_i \otimes (V\mathbf{h}_i))^T \text{vec}(U). \end{aligned} \quad (3.5)$$

Then, we define

$$\mathbf{p}_i = U\mathbf{w}_i, \mathbf{q}_i = V\mathbf{h}_i, \tilde{y}_i = \mathbf{p}_i^T \mathbf{q}_i, b_i = \tilde{y}_i - y_i, \quad (3.6)$$

the loss function

$$\xi_i = b_i^2 = \left( \mathbf{p}_i^T \mathbf{q}_i - y_i \right)^2,$$

$$\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_l] \in \mathcal{R}^{d \times l}, \mathbf{P} = [\mathbf{p}_1, \dots, \mathbf{p}_l] \in \mathcal{R}^{d \times l},$$

$$\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_l]^T \in \mathcal{R}^{l \times M}, \mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_l]^T \in \mathcal{R}^{l \times N}. \quad (3.7)$$

### 3.1 Gradient Calculation

The gradient of (1.2) is

$$\nabla f(\text{vec}(U, V)) = \lambda \text{vec}(U, V) + \frac{1}{2} \sum_{i=1}^l \frac{\partial \xi_i}{\partial \tilde{y}_i} \begin{bmatrix} \frac{\partial \tilde{y}_i}{\partial \text{vec}(U)} \\ \frac{\partial \tilde{y}_i}{\partial \text{vec}(V)} \end{bmatrix}. \quad (3.8)$$

From (3.5), (3.4) and (3.6), we define the following Jacobian of  $\tilde{y}_i$  as

$$\begin{aligned} \mathbf{j}_i &= \begin{bmatrix} \frac{\partial \tilde{y}_i}{\partial \text{vec}(U)} \\ \frac{\partial \tilde{y}_i}{\partial \text{vec}(V)} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\partial (\mathbf{w}_i \otimes (V \mathbf{h}_i))^T \text{vec}(U)}{\partial \text{vec}(U)} \\ \frac{\partial (\mathbf{h}_i \otimes (U \mathbf{w}_i))^T \text{vec}(V)}{\partial \text{vec}(V)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{w}_i \otimes \mathbf{q}_i \\ \mathbf{h}_i \otimes \mathbf{p}_i \end{bmatrix}. \end{aligned} \quad (3.9)$$

From (3.8) and (3.9), the gradient is

$$\begin{aligned}\nabla f(\text{vec}(U, V)) &= \lambda \text{vec}(U, V) + \sum_{i=1}^l b_i \mathbf{j}_i \\ &= \lambda \text{vec}(U, V) + \begin{bmatrix} \sum_{i=1}^l b_i (\mathbf{w}_i \otimes \mathbf{q}_i) \\ \sum_{i=1}^l b_i (\mathbf{h}_i \otimes \mathbf{p}_i) \end{bmatrix}\end{aligned}\quad (3.10)$$

$$= \lambda \text{vec}(U, V) + \begin{bmatrix} \sum_{i=1}^l \text{vec}(\mathbf{q}_i b_i \mathbf{w}_i^T) \\ \sum_{i=1}^l \text{vec}(\mathbf{p}_i b_i \mathbf{h}_i^T) \end{bmatrix}\quad (3.11)$$

$$\begin{aligned}&= \lambda \text{vec}(U, V) + \begin{bmatrix} \text{vec}\left(\sum_{i=1}^l \mathbf{q}_i b_i \mathbf{w}_i^T\right) \\ \text{vec}\left(\sum_{i=1}^l \mathbf{p}_i b_i \mathbf{h}_i^T\right) \end{bmatrix} \\ &= \lambda \text{vec}(U, V) + \begin{bmatrix} \text{vec}\left(Q(\text{diag}(\mathbf{b})W)\right) \\ \text{vec}\left(P(\text{diag}(\mathbf{b})H)\right) \end{bmatrix},\end{aligned}\quad (3.12)$$

where  $\text{diag}(\mathbf{b})$  is a diagonal matrix in which the  $i$ th diagonal element is  $b_i$ , and for the second term in (3.11) we apply (3.1). Usually  $W$  and  $H$  is highly sparse matrices. The cost of gradient calculation is

$$\mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H))),$$

where  $\text{nnz}(\cdot)$  is the number of non-zero elements in a sparse matrix.

### 3.2 Hessian and Gauss-Newton Matrices



The Hessian matrix of (1.2) is

$$\begin{aligned}
 & \frac{\partial}{\partial \text{vec}(U, V)^T} \frac{\partial f}{\partial \text{vec}(U, V)} \\
 &= \lambda I + \sum_{i=1}^l \frac{\partial}{\partial \text{vec}(U, V)^T} (b_i \mathbf{j}_i) \\
 &= \lambda I + \sum_{i=1}^l \left( \frac{\partial}{\partial b_i^T} (b_i \mathbf{j}_i) \frac{\partial b_i}{\partial \text{vec}(U, V)^T} + \frac{\partial}{\partial \mathbf{j}_i^T} (b_i \mathbf{j}_i) \frac{\partial \mathbf{j}_i}{\partial \text{vec}(U, V)^T} \right) \quad (3.13) \\
 &= \lambda I + \sum_{i=1}^l \left( \mathbf{j}_i \frac{\partial b_i}{\partial \tilde{y}_i} \frac{\partial \tilde{y}_i}{\partial \text{vec}(U, V)^T} + b_i I \frac{\partial \mathbf{j}_i}{\partial \text{vec}(U, V)^T} \right) \\
 &= \lambda I + \sum_{i=1}^l \mathbf{j}_i \mathbf{j}_i^T + \sum_{i=1}^l b_i \frac{\partial \mathbf{j}_i}{\partial \text{vec}(U, V)^T},
 \end{aligned}$$

where  $I$  is the identity matrix. This matrix may not be positive definite. We remove the last term in (3.13) and obtain the following Gauss-Newton matrix, which is positive definite.

$$G = \lambda I + \sum_{i=1}^l \mathbf{j}_i \mathbf{j}_i^T \quad (3.14)$$

To get the product between the Gauss-Newton matrix and a vector

$$\mathbf{s} = \text{vec}(S) = \begin{bmatrix} \text{vec}(S_u) \\ \text{vec}(S_v) \end{bmatrix},$$

from (3.14), (3.9) and (3.1), we first calculate

$$\begin{aligned}
 z_i &= \mathbf{j}_i^T \text{vec}(S) \\
 &= \begin{bmatrix} \mathbf{w}_i^T \otimes \mathbf{q}_i^T & \mathbf{h}_i^T \otimes \mathbf{p}_i^T \end{bmatrix} \begin{bmatrix} \text{vec}(S_u) \\ \text{vec}(S_v) \end{bmatrix} \quad (3.15) \\
 &= \mathbf{q}_i^T S_u \mathbf{w}_i + \mathbf{p}_i^T S_v \mathbf{h}_i, \quad i = 1, \dots, l,
 \end{aligned}$$

or equivalently

$$\mathbf{z} = \left( Q^T \odot (W S_u^T) + P^T \odot (H S_v^T) \right) \mathbf{1}_{d \times 1}, \quad (3.16)$$

where  $\odot$  is the Hadamard product (i.e., element-wise product) of two matrices. Next, from (3.15), (3.14), (3.9), (2.4) and (3.1),

$$\begin{aligned}
G\mathbf{s} &= \lambda \text{vec}(S) + \sum_{i=1}^l z_i \mathbf{j}_i \\
&= \lambda \text{vec}(S) + \sum_{i=1}^l z_i \begin{bmatrix} \mathbf{w}_i \otimes \mathbf{q}_i \\ \mathbf{h}_i \otimes \mathbf{p}_i \end{bmatrix} \\
&= \lambda \text{vec}(S) + \sum_{i=1}^l \begin{bmatrix} \text{vec}(z_i \mathbf{q}_i \mathbf{w}_i^T) \\ \text{vec}(z_i \mathbf{p}_i \mathbf{h}_i^T) \end{bmatrix} \\
&= \lambda \text{vec}(S) + \begin{bmatrix} \text{vec}(Q(\text{diag}(\mathbf{z})W)) \\ \text{vec}(P(\text{diag}(\mathbf{z})H)) \end{bmatrix}. \tag{3.17}
\end{aligned}$$

We investigate the complexity of the Hessian-vector product. The first term of (3.17) is a standard vector scaling that costs  $\mathcal{O}((M+N) \times d)$ . For the second term in (3.17), we separately consider two major steps:

1. The computation of  $WS_u^T$  and  $HS_v^T$ :  $\mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H)))$
2. The element-wise product between  $P^T$  and  $HS_v^T$ :  $\mathcal{O}(d \times l)$
3. The product between dense vector  $Q$  and  $\text{diag}(\mathbf{z})W$ :  $\mathcal{O}(d \times l)$

Therefore, the cost of one Hessian-vector is  $\mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H)))$ .

### 3.3 Implementation of Line Search

Recalculating the function value at each  $U + \theta S_u$  and  $V + \theta S_v$  is expensive, where the main cost is on calculating  $((U + \theta S_u)\mathbf{w}_i)^T((V + \theta S_v)\mathbf{h}_i)$ ,  $\forall i$ . However, the following trick in Yuan et al., (2010) can be employed. Assume that

$$\tilde{y}_i = (U\mathbf{w}_i)^T(V\mathbf{h}_i), \tag{3.18}$$

$$\Delta_i = (S_u \mathbf{w}_i)^T (V \mathbf{h}_i) + (U \mathbf{w}_i)^T (S_v \mathbf{h}_i) \quad (3.19)$$

and

$$\Delta'_i = (S_u \mathbf{w}_i)^T (S_v \mathbf{h}_i), \quad i = 1, \dots, l \quad (3.20)$$

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

$$((U + \theta S_u) \mathbf{w}_i)^T ((V + \theta S_v) \mathbf{h}_i) = \tilde{y}_i + \theta \Delta_i + \theta^2 \Delta'_i \quad (3.21)$$

to get the new output value. Now we have

$$\begin{aligned} f(U + \theta S_u, V + \theta S_v) - f(U, V) &= \frac{\lambda}{2} (\|U + S_u\|_F^2 + \|V + S_v\|_F^2 - \|U\|_F^2 - \|V\|_F^2) \\ &\quad + \frac{1}{2} \sum_{i=1}^l (y_i - \tilde{y}_i - \theta \Delta_i - \theta^2 \Delta'_i)^2 - \frac{1}{2} \sum_{i=1}^l (y_i - \tilde{y}_i)^2 \\ &= \frac{\lambda}{2} (2\theta \langle U, S_u \rangle + 2\theta \langle V, S_v \rangle + \theta^2 \|S\|_F^2) \\ &\quad + \frac{1}{2} \sum_{i=1}^l (y_i - \tilde{y}_i - \theta \Delta_i - \theta^2 \Delta'_i)^2 - \frac{1}{2} \sum_{i=1}^l (y_i - \tilde{y}_i)^2, \end{aligned}$$

where  $\langle \cdot, \cdot \rangle$  is the inner product of two matrices. If we further maintain

$$\langle U, S_u \rangle, \langle V, S_v \rangle, \|S\|_F^2, \text{ and } \mathbf{g}^T \mathbf{s}, \quad (3.22)$$

then the total cost of checking the condition in (2.6) is merely  $\mathcal{O}(l)$ .

### 3.4 Overall Procedure

The overall procedure is in Algorithm 1. For the sack of efficiency, if possible, for each operation, we use a matrix-based form rather than a vector-based representation. We discuss some important ones.

For the gradient, from the vector form in (3.12), we rewrite it as a matrix in Line 7. For the stopping condition, we follow past works on Newton methods such as Lin et al., (2008) to use

$$\|\nabla f(\text{vec}(U), \text{vec}(V))\| \leq \epsilon \|\nabla f(\text{vec}(U_{\text{init}}), \text{vec}(V_{\text{init}}))\|, \quad (3.23)$$



where  $0 < \epsilon < 1$  is a small stopping tolerance and  $(U_{\text{init}}, V_{\text{init}})$  indicates the initial point of the model. By the matrix form in Line 7, this stopping condition is essentially

$$\|\nabla f(U, V)\|_F \leq \epsilon \|\nabla f(U_{\text{init}}, V_{\text{init}})\|_F, \quad (3.24)$$

where  $\|\cdot\|_F$  is Frobenius norm of a matrix. For  $\Delta_i, \Delta'_i$  in (3.21), for line search, we aggregate all values together to have the form in Line 16. Besides, other values in (3.22) to be maintained are in Line 17. The check of the sufficient decrease condition (2.6) is by the calculation in Line 19.

For details of the CG procedure to approximately solve the linear system (2.3), we provide Algorithm 2. The main cost in each iteration of CG is the product between the Gauss-Newton matrix and a vector in Line 4 and 5, where we follow details in (3.16) and (3.17). Its computational complexity is  $\mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H)))$  from what we discussed in (3.17).

After  $s_k$  is obtained from the conjugate gradient method, a backtracking line search is applied to adjust the step size  $\alpha$ . From the discussion in Section 3.3, the complexity of the line search is

$$\mathcal{O}((\# \text{ of line search steps}) \times l).$$

The overall complexity per Newton iteration is

$$(\# \text{ of CG iterations} + 2) \times \mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H))) + (\# \text{ of line search steps}) \times \mathcal{O}(l), \quad (3.25)$$

where the term  $2 \times \mathcal{O}(d \times (\text{nnz}(W) + \text{nnz}(H)))$  is from the cost in Line 7 and Line 16 in Algorithm 1 to compute (3.12), (3.19) and (3.20). If matrix factorization is considered, (1.4) and (1.5) imply that

$$\text{nnz}(W) = \text{nnz}(H) = l. \quad (3.26)$$

Thus (3.25) can be written as

$$(\# \text{ of CG iterations} + 2) \times \mathcal{O}(ld) + (\# \text{ of line search steps}) \times \mathcal{O}(l).$$

Next we discuss the memory usage. The main bottleneck occurs when conducting Line 4 of Algorithm 2 or Line 16 of Algorithm 1. Take Line 4 of Algorithm 2 as an example. We must store the following matrix

$$P \in \mathcal{R}^{d \times l}, Q \in \mathcal{R}^{d \times l}, W D_u^T \in \mathcal{R}^{l \times d}, H D_v^T \in \mathcal{R}^{l \times d}, W \in \mathcal{R}^{l \times M}, \text{ and } H \in \mathcal{R}^{l \times N}. \quad (3.27)$$

Thus the total memory consumption is

$$4 \times \mathcal{O}(ld) + \text{nnz}(W) + \text{nnz}(H).$$

For the case of matrix factorization, from (3.26),  $\text{nnz}(W)$  and  $\text{nnz}(H)$  are smaller, so the four  $l \times d$  matrices are the bottleneck.

### 3.5 Implementation for Matrix Factorization

From (1.4), (1.5) and (3.7), for matrix factorization  $W$  and  $H$  are special matrices when each row contains exactly one non-zero entry. We explore how to take such properties to reduce the memory consumption.

Consider the two main operations (3.16) and (3.17) for the Gauss-Newton matrix-vector product. From

$$P = U W^T \text{ and } Q = V H^T,$$

(3.16) is the same as

$$\mathbf{z} = \left( (H V^T) \odot (W S_u^T) + (W U^T) \odot (H S_v^T) \right) \mathbf{1}_{d \times 1}.$$

Let  $\mathbf{z}$  be indexed by  $(m, n)$  in  $R$  with  $R_{m,n} \neq 0$ , and

$$S_u = [\mathbf{S}_u^1, \dots, \mathbf{S}_u^M] \in \mathcal{R}^{d \times M},$$

$$S_v = [\mathbf{S}_v^1, \dots, \mathbf{S}_v^N] \in \mathcal{R}^{d \times N}.$$

Then each component in (3.16) is

$$\mathbf{z}_{(m,n)} = \mathbf{v}_n^T \mathbf{S}_u^m + \mathbf{u}_m^T \mathbf{S}_v^n.$$

Let  $Z \in \mathcal{R}^{M \times N}$  be a sparse matrix to include these  $\mathbf{z}_{(m,n)}$  elements.

We note that

$$H^T W = I [R^T],$$

where  $I [\cdot]$  is an indicator function with values 1 or 0 to reflect non-zero position of the input matrix. Thus in the calculation of (3.17) we have

$$\begin{aligned} & Q \text{diag}(\mathbf{z}) W \\ &= V H^T \text{diag}(\mathbf{z}) W \\ &= V (Z^T \odot I [R^T]). \end{aligned}$$

Further

$$\begin{aligned} & P \text{diag}(\mathbf{z}) H \\ &= U W^T \text{diag}(\mathbf{z}) H \\ &= U (Z \odot I [R]). \end{aligned}$$

By the above settings, we never need to store any matrix in (3.27). Therefore the memory consumption is significantly reduced to

$$\mathcal{O}(\text{nnz}(R) + (M + N)d).$$

For the computational complexity, it remains the same as the one shown in (3.25). The reason is that

$$\text{nnz}(W) = \text{nnz}(H) = \text{nnz}(R).$$

We give details of other places in Algorithm 2 that also need to be changed. Let  $B \in \mathcal{R}^{M \times N}$  contains elements  $B_{m,n} = \mathbf{u}_m^T \mathbf{v}_n - r_{m,n}$ . At Line 7 in Algorithm 1, we have

$$G \leftarrow \lambda \left[ U V \right] + \left[ V (B^T \odot I [R^T]) U (B \odot I [R]) \right].$$



Then let  $Z' \in \mathcal{R}^{M \times N}$  contains elements  $Z'_{m,n} = (\mathbf{S}_u^m)^T \mathbf{S}_v^n$ . At Line 16 in Algorithm 1, we have

$$\Delta = Z,$$

$$\Delta' = Z'.$$






---

**Algorithm 1** Newton method to solve (1.2) by matrix-based operations.

---

```

1: Given  $0 < \epsilon < 1$  and the initial  $(U, V)$ .
2: Compute and cache  $\tilde{\mathbf{y}} = ((VH^T) \odot (UW^T))^T \mathbf{1}_{d \times 1}$ .
3:  $\mathbf{b} \leftarrow \tilde{\mathbf{y}} - \mathbf{y}$ 
4:  $f \leftarrow \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \mathbf{b}^T \mathbf{b}$ .
5: for  $k \leftarrow \{0, 1, \dots\}$  do
6:    $Q \leftarrow VH^T, P \leftarrow UW^T$ 
7:   Compute  $G \leftarrow \lambda [UV] + [Q(\text{diag}(\mathbf{b})W) P(\text{diag}(\mathbf{b})H)]$ .
8:   if  $k = 0$  then
9:      $\|G^0\|_F \leftarrow \|G\|_F$ 
10:  end if
11:  if  $\|G\|_F \leq \epsilon \|G^0\|_F$  then
12:    Output  $[UV]$  as the solution of (1.2).
13:    break
14:  end if
15:  Run CG in Algorithm 5 to get an update direction  $S = [S_u S_v]$ .
16:  Calculate

```

$$\Delta = (Q^T \odot (WS_u^T) + P^T \odot (HS_v^T)) \mathbf{1}_{d \times 1},$$

$$\Delta' = ((WS_u^T) \odot (HS_v^T)) \mathbf{1}_{d \times 1}.$$

```

17:   Parpare the following values

```

$$\langle U, S_u \rangle, \langle V, S_v \rangle, \|S\|_F^2, \text{ and } \langle G, S \rangle.$$

```

18:   for  $\theta \leftarrow \{1, \beta, \beta^2, \dots\}$  do
19:      $\delta \leftarrow \frac{\lambda}{2} (2\theta \langle U, S_u \rangle + 2\theta \langle V, S_v \rangle + \theta^2 \|S\|_F^2) + \frac{1}{2} \|\mathbf{b} - \theta \Delta - \theta^2 \Delta'\|^2 - \frac{1}{2} \mathbf{b}^T \mathbf{b}$ 
20:     if  $\delta \leq \theta \nu \langle G, S \rangle$  then
21:        $U \leftarrow U + \theta S_u, V \leftarrow V + \theta S_v$ 
22:        $f \leftarrow f + \delta$ 
23:        $\tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{y}} + \theta \Delta + \theta^2 \Delta'$ 
24:        $\mathbf{b} \leftarrow \tilde{\mathbf{y}} - \mathbf{y}$ 
25:       break
26:     end if
27:   end for
28: end for

```

---




---

**Algorithm 2** A conjugate gradient method for solving (2.3) by matrix-based operations.

---

- 1: Given  $0 < \eta < 1$  and  $G$ , the gradient matrix form of (3.12). Let  $S = \mathbf{0}_{d \times (M+N)}$ .
  - 2: Calculate  $R = -G$ ,  $D = R$ , and  $\gamma^0 = \gamma = \|R\|_F^2$ .
  - 3: **while**  $\sqrt{\gamma} > \eta\sqrt{\gamma^0}$  **do**
  - 4:    $\mathbf{z} \leftarrow (Q^T \odot (WD_u^T) + P^T \odot (HD_v^T)) \mathbf{1}_{d \times 1}$
  - 5:    $D_h \leftarrow (\lambda D + [Q(\text{diag}(\mathbf{z})W) P(\text{diag}(\mathbf{z})H)])$
  - 6:    $\alpha \leftarrow \gamma / \langle D, D_h \rangle$
  - 7:    $S \leftarrow S + \alpha D$
  - 8:    $R \leftarrow R - \alpha D_h$
  - 9:    $\gamma^{\text{new}} \leftarrow \|R\|_F^2$
  - 10:    $\beta \leftarrow \gamma^{\text{new}} / \gamma$
  - 11:    $D \leftarrow R + \beta D$
  - 12:    $\gamma \leftarrow \gamma^{\text{new}}$
  - 13: **end while**
  - 14: Output  $S$  as the solution.
-



## CHAPTER IV

### A Review of Alternating Newton Methods

In Chin et al., (2018), they solve the optimization problem (1.2) by a block coordinate descent method. For the two blocks  $U$  and  $V$ , each time one block is fixed while the other is updated. A simple illustration of the procedure is in Algorithm 3.

For each sub-problem to update one block, Chin et al., (2018) apply a truncated Newton method. Note that because (1.3) is multi-block convex function, each sub-problem is convex. We consider the sub-problem of  $U$  as an example to show details of the Newton method. Instead of copying results from Chin et al., (2018), here we show that materials can be extracted from the more sophisticated algorithm in Section III for minimizing both blocks together. To begin, from (3.8)-(3.12), we have

$$\nabla_U f(U, V) = \lambda U + Q(\text{diag}(\mathbf{b})W). \quad (4.1)$$

For the Hessian matrix, from (3.13),

$$\frac{\partial}{\partial \text{vec}(U)^T} \frac{\partial f}{\partial \text{vec}(U)} = \lambda I + \sum_{i=1}^l (\mathbf{w}_i \otimes \mathbf{q}_i)(\mathbf{w}_i \otimes \mathbf{q}_i)^T + \sum_{i=1}^l b_i \frac{\partial(\mathbf{w}_i \otimes \mathbf{q}_i)}{\partial \text{vec}(U)^T} \quad (4.2)$$

$$= \lambda I + \sum_{i=1}^l (\mathbf{w}_i \otimes \mathbf{q}_i)(\mathbf{w}_i \otimes \mathbf{q}_i)^T. \quad (4.3)$$

Note that the last term in (4.2) is zero because  $\mathbf{w}_i \otimes \mathbf{q}_i$  is not a function of  $U$ . Therefore, the Hessian matrix is positive definite. Further, it is a constant matrix independent of  $U$ . The

---

**Algorithm 3** Solving problem (1.2) via alternating minimization.

---

```

1: Given an initial solution  $(U, V)$ 
2: while stopping condition is not satisfied do
3:    $U \leftarrow \arg \min_U f(U, V)$ 
4:    $V \leftarrow \arg \min_V f(U, V)$ 
5: end while

```

---



reason is that when  $V$  is fixed, the square loss function considered in (1.3) leads to a convex quadratic function of  $U$ .

The sub-problem of  $U$  becomes equivalent to solving the following linear system by the conjugate gradient method.

$$\left( \frac{\partial}{\partial \text{vec}(U)^T} \frac{\partial f}{\partial \text{vec}(U)} \right) \mathbf{s} = - \frac{\partial f}{\partial \text{vec}(U)}. \quad (4.4)$$

To get the product between the Hessian matrix and a vector

$$\mathbf{s} = \text{vec}(S_u), \quad (4.5)$$

from (4.3), (3.1), and (3.15), we first calculate

$$\begin{aligned} z_i &= \left[ \mathbf{w}_i^T \otimes \mathbf{q}_i^T \right] \text{vec}(S_u) \\ &= \mathbf{q}_i^T S_u \mathbf{w}_i, \quad i = 1, \dots, l, \end{aligned} \quad (4.6)$$

or equivalently

$$\mathbf{z} = \left( Q^T \odot (W S_u^T) \right) \mathbf{1}_{d \times 1}. \quad (4.7)$$

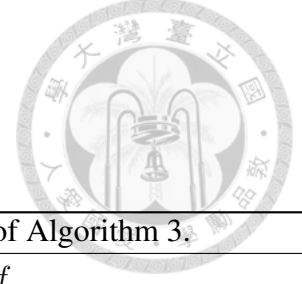
Let  $H$  be the Hessian in (4.3). From (4.6), (4.3), (4.5) and (3.1),

$$\begin{aligned} H \mathbf{s} &= \lambda \text{vec}(S_u) + \sum_{i=1}^l z_i (\mathbf{w}_i \otimes \mathbf{q}_i) \\ &= \lambda \text{vec}(S_u) + \sum_{i=1}^l \text{vec} (z_i \mathbf{q}_i \mathbf{w}_i^T) \\ &= \lambda \text{vec}(S_u) + \text{vec} \left( Q (\text{diag}(\mathbf{z}) W) \right). \end{aligned} \quad (4.8)$$

The computational complexity is

$$(\# \text{ of CG iterations} + 2) \times \mathcal{O}(d \times \text{nnz}(W)).$$






---

**Algorithm 4** Newton method for solving a sub-problem in step 3 of Algorithm 3.

---

```

1: Given  $0 < \epsilon < 1$ ,  $0 < \eta < 1$ ,  $\mathbf{y}$ ,  $W$ ,  $H$ , and the current  $U$ ,  $V$ ,  $f$ .
2:  $Q \leftarrow V H^T$ 
3: Compute and cache  $\tilde{\mathbf{y}} = (Q \odot U W^T)^T \mathbf{1}_{d \times 1}$ .
4:  $\mathbf{b} \leftarrow \tilde{\mathbf{y}} - \mathbf{y}$ 
5: for  $k \leftarrow \{0, 1, \dots\}$  do
6:   Compute  $G \leftarrow \lambda U + Q(\text{diag}(\mathbf{b})W)$ 
7:   if  $k = 0$  then
8:      $\|G^0\|_F \leftarrow \|G\|_F$ 
9:   end if
10:  if  $\|G\|_F \leq \epsilon \|G^0\|_F$  then
11:    Output  $U$  as the solution of step 3 in Algorithm 3.
12:    break
13:  end if
14:  Let  $S_u = \mathbf{0}_{d \times M}$ .
15:  Calculate  $R = -G$ ,  $D = R$ , and  $\gamma^0 = \gamma = \|R\|_F^2$ .
16:  while  $\sqrt{\gamma} > \eta \sqrt{\gamma^0}$  do
17:     $\mathbf{z} \leftarrow (Q^T \odot (W D^T)) \mathbf{1}_{d \times 1}$ 
18:     $D_h \leftarrow (\lambda D + Q(\text{diag}(\mathbf{z})W))$ 
19:     $\alpha \leftarrow \gamma / \langle D, D_h \rangle$ 
20:     $S_u \leftarrow S_u + \alpha D$ 
21:     $R \leftarrow R - \alpha D_h$ 
22:     $\gamma^{\text{new}} \leftarrow \|R\|_F^2$ 
23:     $\beta \leftarrow \gamma^{\text{new}} / \gamma$ 
24:     $D \leftarrow R + \beta D$ 
25:     $\gamma \leftarrow \gamma^{\text{new}}$ 
26:  end while
27:  Calculate  $\Delta = (Q^T \odot (W S_u^T)) \mathbf{1}_{d \times 1}$ 
28:   $\delta \leftarrow \frac{\lambda}{2} (2 \langle U, S_u \rangle + \|S_u\|_F^2) + \frac{1}{2} \|\mathbf{b} - \theta \Delta\|^2 - \frac{1}{2} \mathbf{b}^T \mathbf{b}$ 
29:   $U \leftarrow U + S_u$ 
30:   $f \leftarrow f + \delta$ 
31:   $\tilde{\mathbf{y}} \leftarrow \tilde{\mathbf{y}} + \theta \Delta$ 
32:   $\mathbf{b} \leftarrow \tilde{\mathbf{y}} - \mathbf{y}$ 
33: end for

```

---

The CG procedure stops if the CG iterate  $\mathbf{s}$  satisfies

$$\| H\mathbf{s} + \frac{\partial f}{\partial \text{vec}(U)} \| \leq \eta \| \frac{\nabla f}{\partial \text{vec}(U)} \|. \quad (4.9)$$

Note that because the sub-problem is equivalent to a linear system (4.4), one single CG procedure is enough to obtain an approximate solution. Thus there is no need to have an outer Newton procedure. The CG procedure is presented in Algorithm 4. Clearly, these procedures are simplified from algorithms considering  $U$  and  $V$  together. Note that in Algorithm 4, we do not need to form  $P \leftarrow UW^T$  as in Algorithm 1 because  $UW^T$  is used only once.

To solve the sub-problem over  $V$ , we apply the same procedure through the following values which are swapped.

$$U \leftrightarrow V$$

$$S_u \leftrightarrow S_v$$

$$\mathbf{w}_i \leftrightarrow \mathbf{h}_i$$

$$\mathbf{q}_i \leftrightarrow \mathbf{p}_i$$

$$W \leftrightarrow H$$

$$Q \leftrightarrow P$$

Regarding the memory cost, the bottleneck is at similar places to those in Algorithm 1. However, the needed memory is only about half of Algorithm 1. This can be clearly seen from comparing Line 4 of Algorithm 2 and Line 17 of Algorithm 4.



## CHAPTER V

### Experiments

We conduct experiments to compare the following two methods to solve (1.2).

- Newton method described in Section II.
- Alternating Newton method (Chin et al., 2018) described in Section III.

#### 5.1 Experimental Settings

Table 5.1: Statistics and parameters for each data set.

Data Set	movielens10m	netflix	movielens1m
$M$	71,567	2,649,429	6,040
$N$	65,133	17,770	3,952
#Training	9,301,274	99,072,112	*
#Test	698,780	1,408,395	*
$d$	40	20	40

We consider data sets listed in Table 5.1. For both methods, every element of initial  $U$  and  $V$  is randomly drawn from the interval  $[-0.1/\sqrt{d}, 0.1/\sqrt{d}]$ .

For the CG procedure, we set  $\eta = 0.3$  as the stopping tolerance. Note that we consider the same  $\eta$  for the CG procedure in both Algorithms 2 and ???. For the line search procedure, we set  $\beta = 0.5$  as the ratio to decrease the step size and  $\nu = 0.1$  for the sufficient decrease condition (2.6). For the latent dimension, we use  $d$  listed in Table 5.1.

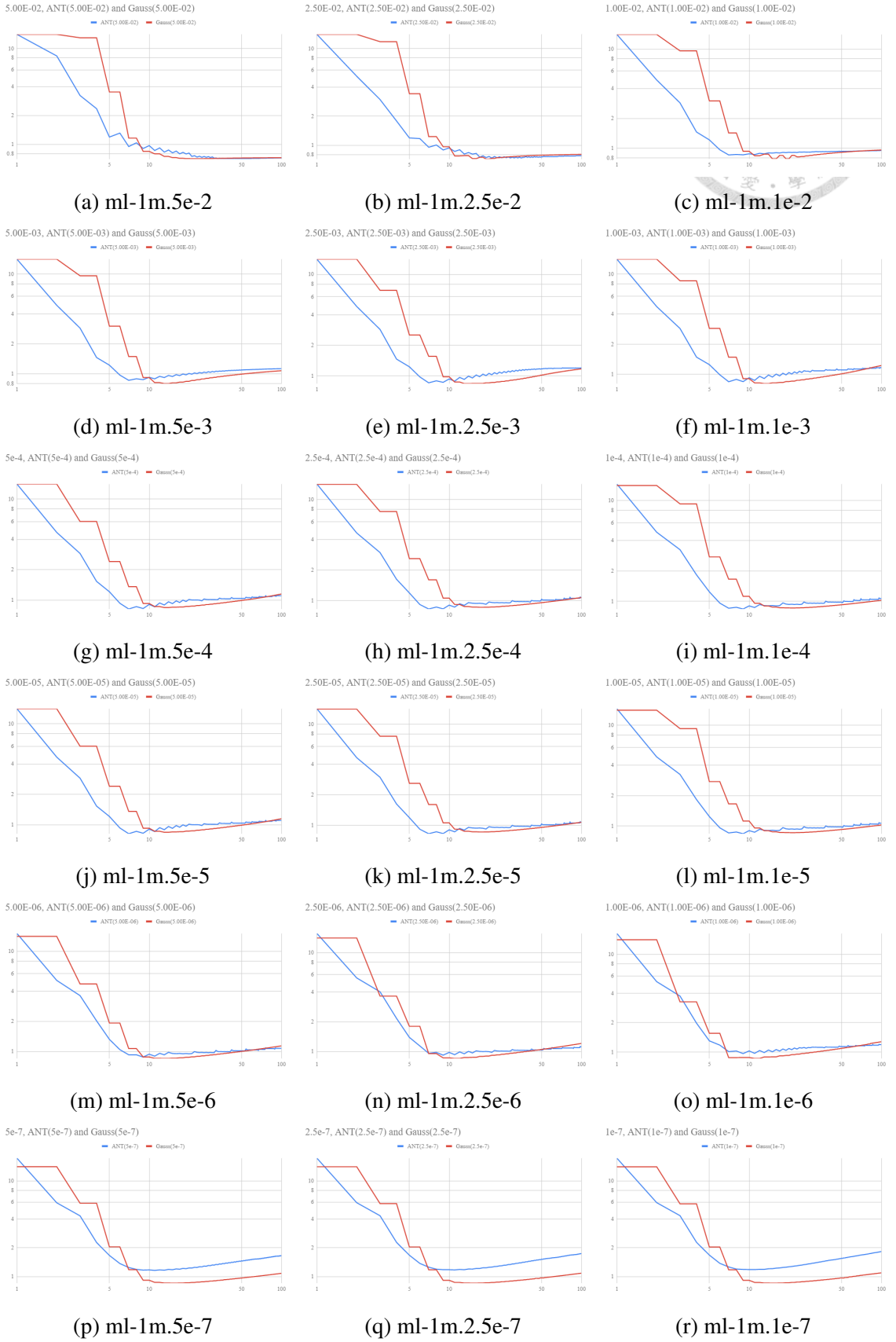


Figure 5.1: A comparison on the va-loss of different methods. ANT is worse than Gauss among smaller  $\lambda$ . The  $x$ -axis is the iteration, while the  $y$ -axis is the va-loss.

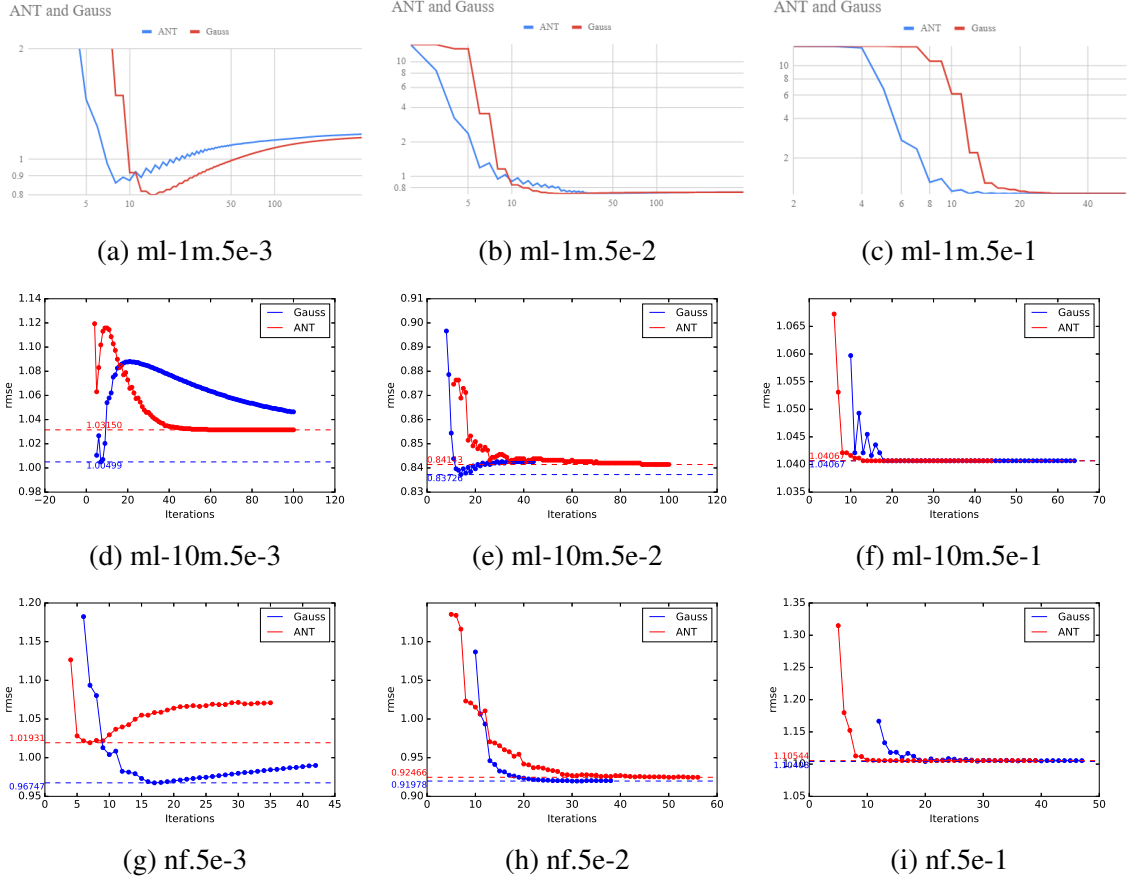
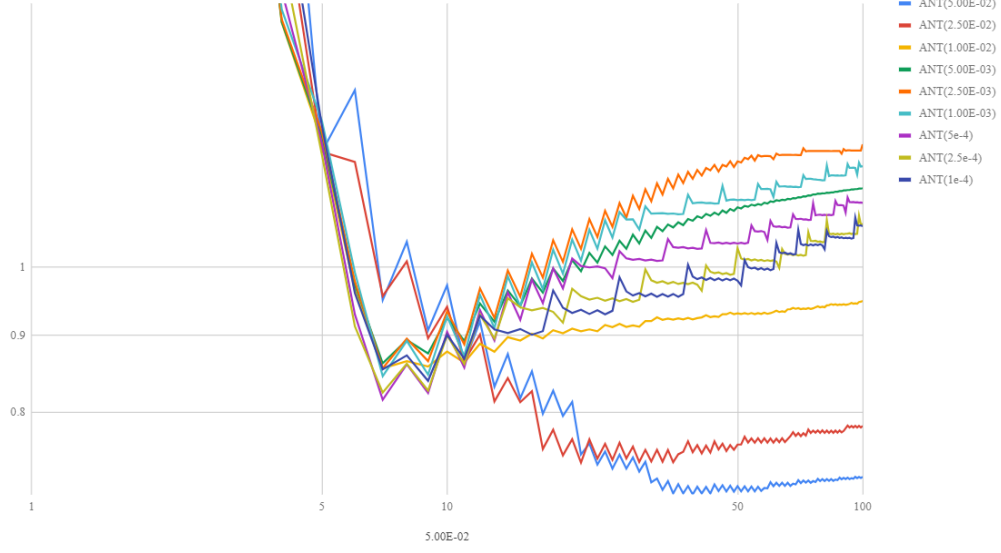


Figure 5.2: A comparison on the va-loss of different methods. The  $x$ -axis is the iteration, while the  $y$ -axis is the va-loss.

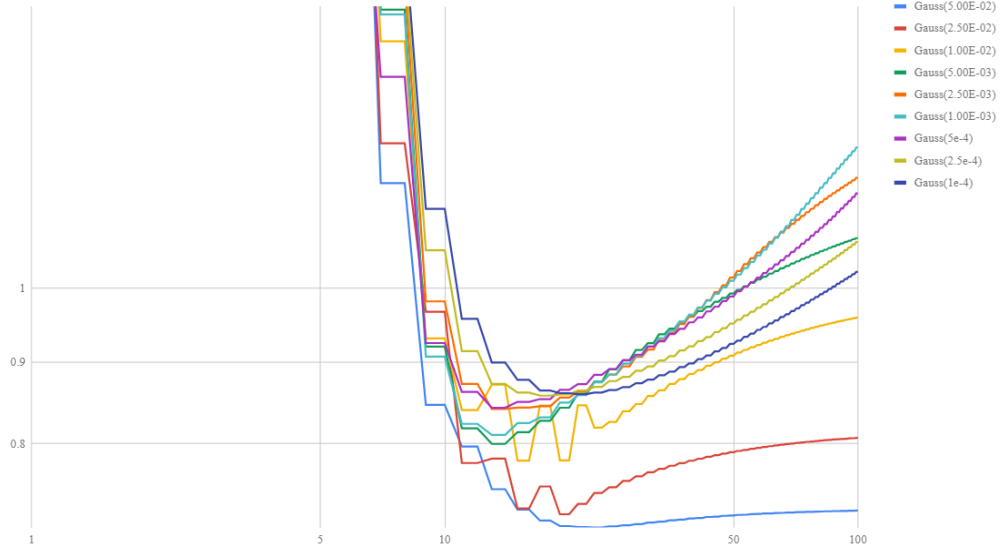


ANT



(a) Ant

Gauss



(b) Gauss

Figure 5.3: ml-1m. A comparison on the va-loss of different  $\lambda$ . The  $x$ -axis is the iteration, while the  $y$ -axis is the va-loss.



Figure 5.4: ml-1m. A comparison on the best va-loss of different methods. The  $x$ -axis is  $\lambda$ , while the  $y$ -axis is the va-loss.



## Bibliography

- Rendle, Steffen (2010). "Factorization machines". In: *Proceedings of IEEE International Conference on Data Mining (ICDM)*, pp. 995–1000.
- Blondel, Mathieu et al. (2016). "Polynomial Networks and Factorization Machines: new Insights and Efficient Training Algorithms". In: *Proceedings of the 33rd International Conference on Machine Learning (ICML)*.
- Chin, Wei-Sheng et al. (2018). "An Efficient Alternating Newton Method for Learning Factorization Machines". In: *ACM Transactions on Intelligent Systems and Technology* 9.6, 72:1–72:31. URL: <https://www.csie.ntu.edu.tw/~cjlin/papers/fm/scalefm.pdf>.
- Yuan, Guo-Xun et al. (2010). "A Comparison of Optimization Methods and software for Large-scale L1-regularized Linear Classification". In: *Journal of Machine Learning Research* 11, pp. 3183–3234. URL: <http://www.csie.ntu.edu.tw/~cjlin/papers/l1.pdf>.
- Lin, Chih-Jen, Ruby C. Weng, and S. Sathiya Keerthi (2008). "Trust region Newton method for large-scale logistic regression". In: *Journal of Machine Learning Research* 9, pp. 627–650. URL: <http://www.csie.ntu.edu.tw/~cjlin/papers/logistic.pdf>.



We consider a preconditioner  $\tilde{M}$  to approximately factorize  $H$  such that  $H \approx \tilde{M}\tilde{M}^T$ , and then use CG to solve

$$\hat{H}\hat{s} = \hat{g}, \quad (5.1)$$

where  $\hat{H} = \tilde{M}^{-1}H\tilde{M}^{-T}$  and  $\hat{g} = \tilde{M}^{-1}g$ . Once  $\hat{s}$  is found, the solution of (2.2) can be recovered by  $s = \tilde{M}^{-T}\hat{s}$ . Then, we consider the following diagonal preconditioner for solving (2.2) as an example.

$$\tilde{M} = \tilde{M}^T = \text{diag}(\mathbf{h}), \quad (5.2)$$

where

$$\mathbf{h} = \sqrt{\lambda \mathbf{1}_{d(M+N) \times 1} + \sum_{i=1}^l \mathbf{j}_i \odot \mathbf{j}_i}.$$

Note that “ $\sqrt{\cdot}$ ” element-wisely performs the square-root operation if the input argument is a vector or a matrix. From

$$\mathbf{j}_i = \begin{bmatrix} \mathbf{w}_i \otimes \mathbf{q}_i \\ \mathbf{h}_i \otimes \mathbf{p}_i \end{bmatrix} = \begin{bmatrix} \text{vec}(\mathbf{q}_i \mathbf{w}_i^T) \\ \text{vec}(\mathbf{p}_i \mathbf{h}_i^T) \end{bmatrix}, \quad (5.3)$$

we have

$$\begin{aligned} \sum_{i=1}^l \mathbf{j}_i \odot \mathbf{j}_i &= \sum_{i=1}^l \begin{bmatrix} \text{vec}((\mathbf{q}_i \odot \mathbf{q}_i)(\mathbf{w}_i \odot \mathbf{w}_i)^T) \\ \text{vec}((\mathbf{p}_i \odot \mathbf{p}_i)(\mathbf{h}_i \odot \mathbf{h}_i)^T) \end{bmatrix} \\ &= \begin{bmatrix} \text{vec}((Q \odot Q)(W \odot W)) \\ \text{vec}((P \odot P)(H \odot H)) \end{bmatrix}. \end{aligned} \quad (5.4)$$

Thus, the preconditioner can be obtained via

$$\tilde{M} = \tilde{M}^T = \text{diag}\left(\sqrt{\lambda \mathbf{1}_{d(M+N) \times 1} + \begin{bmatrix} \text{vec}((Q \odot Q)(W \odot W)) \\ \text{vec}((P \odot P)(H \odot H)) \end{bmatrix}}\right). \quad (5.5)$$




---

**Algorithm 5** A preconditioned conjugate gradient method for solving (2.2) by operations on matrix variables.

---

- 1: Given  $0 < \eta < 1$  and  $G$ , the gradient matrix form of (1.2). Let  $\hat{S} = \mathbf{0}_{d \times (M+N)}$ .
  - 2: Compute  $M$  via (5.5).
  - 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:    $\hat{\mathbf{z}} \leftarrow \left( Q^T \odot (W \hat{D}_{hu}^T) + P^T \odot (H \hat{D}_{hv}^T) \right) \mathbf{1}_{d \times 1}$
  - 7:    $\hat{D}_h \leftarrow \left( \lambda \hat{D}_h + [Q(\text{diag}(\hat{\mathbf{z}})W) P(\text{diag}(\hat{\mathbf{z}})H)] \right) ./ M$
  - 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.
-