

**CME 323: Distributed Algorithms and Optimization, Spring 2015**

<http://stanford.edu/~rezab/dao>.

Instructor: Reza Zadeh, Databricks and Stanford.

**Lecture 14, 5/13/2015. Scribed by Haoming Li, Bangzheng He, Michael Lublin, Yonathan Perez.**

## 14 Matrix Completion via Alternating Least Square(ALS)

### 14.1 Introduction

A common problem faced by internet companies is that of recommending new products to users in personalized settings (e.g. Amazon's product recommender system, and Netflix movie recommendations). This can be formulated as a learning problem in which we are given the ratings that users have given certain items and are tasked with predicting their ratings for the rest of the items. Formally, if there are  $n$  users and  $m$  items, we are given an  $n \times m$  matrix  $R$  in which the  $(u, i)^{th}$  entry is  $r_{ui}$  – the rating for item  $i$  by user  $u$ . Matrix  $R$  has many missing entries indicating unobserved ratings, and our task is to estimate these unobserved ratings.

### 14.2 Matrix Factorization: Objective and ALS Algorithm on a Single Machine

A popular approach for this is **matrix factorization**, where we fix a relatively small number  $k$  (e.g.  $k \approx 10$ ), and summarize each user  $u$  with a  $k$  dimensional vector  $x_u$ , and each item  $i$  with a  $k$  dimensional vector  $y_i$ . These vectors are referred to as *factors*. Then, to predict user  $u$ 's rating for item  $i$ , we simply predict  $r_{ui} \approx x_u^T y_i$ . This can be put in matrix form: Let  $x_1, \dots, x_n \in \mathbb{R}^k$  be the factors for the users, and  $y_1, \dots, y_m \in \mathbb{R}^k$  the factors for the items. The  $k \times n$  user matrix  $X$ , and the  $k \times m$  item matrix  $Y$  are then defined by:

$$X = \begin{bmatrix} | & & | \\ x_1 & \cdots & x_n \\ | & & | \end{bmatrix}, Y = \begin{bmatrix} | & & | \\ y_1 & \cdots & y_m \\ | & & | \end{bmatrix}$$

Our goal is then to estimate the complete ratings matrix  $R \approx X^T Y$ . We can formulate this problem as an optimization problem in which we aim to minimize an objective function and find optimal  $X$  and  $Y$ . In particular, we aim to minimize the least squares error of the observed ratings (and regularize):

$$\min_{X, Y} \sum_{r_{ui} \text{ observed}} (r_{ui} - x_u^T y_i)^2 + \lambda (\sum_u \|x_u\|^2 + \sum_i \|y_i\|^2) \quad (1)$$

Notice that this objective is non-convex (because of the  $x_u^T y_i$  term); in fact it's NP-hard to optimize. Gradient descent can be used as an approximate approach here, however it turns out to be slow and costs lots of iterations. Note however, that if we fix the set of variables  $X$  and treat them as constants, then the objective is a convex function of  $Y$  and vice versa. Our approach will

therefore be to fix  $Y$  and optimize  $X$ , then fix  $X$  and optimize  $Y$ , and repeat until convergence. This approach is known as **ALS(Alternating Least Squares)**. For our objective function, the alternating least squares algorithm is as follows:

---

**Algorithm 1** ALS for Matrix Completion

---

Initialize  $X, Y$

**repeat**

**for**  $u = 1 \dots n$  **do**

$$x_u = \left( \sum_{r_{ui} \in r_{u*}} y_i y_i^T + \lambda I_k \right)^{-1} \sum_{r_{ui} \in r_{u*}} r_{ui} y_i \quad (2)$$

**end for**

**for**  $i = 1 \dots m$  **do**

$$y_i = \left( \sum_{r_{ui} \in r_{*i}} x_u x_u^T + \lambda I_k \right)^{-1} \sum_{r_{ui} \in r_{*i}} r_{ui} x_u \quad (3)$$

**end for**

**until** convergence

---

For a single machine we can analyze the computational cost of this algorithm. Updating each  $x_u$  will cost  $O(n_u k^2 + k^3)$  where  $n_u$  is the number of items rated by user  $u$ , and similarly updating each  $y_i$  will cost  $O(n_i k^2 + k^3)$  where  $n_i$  is the number of users that have rated item  $i$ .

Once we've computed the matrices  $X$  and  $Y$ , there are several ways to do prediction. The first is to do what was discussed before, which is to simply predict  $r_{ui} \approx x_u^T y_i$  for each user  $u$  and item  $i$ . This approach will cost  $O(nmk)$  if we'd like to estimate every user-item pair. However, this approach is prohibitively expensive for most real-world datasets. A second (and more holistic) approach is to use the  $x_u$  and  $y_i$  as features in another learning algorithm, incorporating these features with others that are relevant to the prediction task.

### 14.3 Distributed ALS

There are several ways to distribute the computation of the ALS algorithm(1) depending on how we distribute the data.

#### Method 1 (join)

First we consider a fully distributed version. In this setup all the data (ratings) and parameters ( $X$  and  $Y$ ) are distributed (i.e. stored in an RDD). Specifically, ratings are stored as RDD of triplets:

$$Ratings : RDD((u, i, r_{ui}), \dots)$$

This is reasonable since ratings are always sparse. We use dense representation for factor matrices  $X$  and  $Y$ , and these are stored as RDDs of vectors:

$$X : RDD(x_1, \dots, x_n)$$

$$Y : RDD(y_1, \dots, y_m)$$

Recall equation(2) from the ALS algorithm(1):

$$x_u = \left( \underbrace{\sum_{r_{ui} \in r_{u*}} y_i y_i^T}_A + \lambda I_k \right)^{-1} \underbrace{\sum_{r_{ui} \in r_{u*}} r_{ui} y_i}_B \quad (4)$$

To compute parts A and B, we can follow the steps below:

1. Join *Ratings* with  $Y$  factors using key  $i$  (items)
2. Map to compute  $y_i y_i^T$  and change key to  $u$  (user)
3. ReduceByKey  $u$  (user) to compute  $\sum_i y_i y_i^T$
4. Invert
5. Another ReduceByKey  $u$ (user) to compute  $\sum r_{ui} y_i$

We can use the similar method to compute equation(3) too.

This approach will always work, but note that we are redoing the computation of  $y_i y_i^T$  for each user that has rated item  $i$ . The next approach takes advantage of the fact that the  $X$  and  $Y$  factor matrices are often very small and can be stored locally on each machine:

### Method 2 (broadcast)

1. Partition *Ratings* by user to create  $R_1$ , and similarly partition *Ratings* by item to create  $R_2$  (so there are two copies of *Ratings* with different partitionings). In  $R_1$ , all ratings by the same user are on the same machine, and in  $R_2$  all ratings for same item are on the same machine.
2. Broadcast  $X, Y$  (These are now local matrices and not RDDs)
3. Using  $R_1$  and  $Y$ , we can use equation(2) to compute the update of  $x_u$  locally on each machine
4. Using  $R_2$  and  $X$ , we can use equation(3) to compute the update of  $y_i$  locally on each machine

A further optimization to this method is to group the  $X$  and  $Y$  factors into blocks (user blocks and item blocks) and reduce the communication by only sending to each machine the block of users (or items) that are needed to compute the updates at that machine. This method is called **Block ALS**. It is achieved by precomputing some information about the ratings matrix to determine the “out-links” of each user(which blocks of the items it will contribute to) and “in-link” information for each item(which of the factor vectors it receives from each user block it will depend on). For example, assume that machine 1 is responsible for users 1,2,...,37 – these will be block 1 of users. The items rated by these users are block 1 of items. Only the factors of block 1 of users and block 1 of items will be broadcasted to machine 1.

## 14.4 Streaming ALS

We now assume that ratings data  $r_{ui}$  is streaming in, and that factor matrices  $X, Y$  fit in memory of each machine. We assume that the stream is shuffled and can use Stochastic Gradient Descent(SGD) to update the  $X, Y$  factor matrices.

---

**Algorithm 2** Streaming ALS using SGD

---

```
for new  $r_{ui}$  do  
   $x_u \leftarrow x_u - \alpha(r_{ui} - x_u^\top y_i)y_i + \lambda x_u$   
   $y_i \leftarrow y_i - \alpha(r_{ui} - x_u^\top y_i)x_u + \lambda y_i$   
end for
```

---

## 14.5 Other methods for Matrix Factorization

Other methods for matrix factorization include:

- Sparse Subspace Embedding[1]
- Generalized Low Rank Models(GLRM)[2]
- Fast ALS[3]

## References

- [1] Kenneth L. Clarkson and David P. Woodruff. STOC, 2013.
- [2] Madeleine Udell, Corinne Horn, Reza Zadeh, Stephen Boyd. *Generalized Low Rank Models*.
- [3] Trevor Hastie, Rahul Mazumder, Jason D. Lee, Reza Zadeh *Matrix Completion and Low-Rank SVD via Fast Alternating Least Squares*. Statistics Department and ICME, Stanford University, 2014.