

Advanced Recommender Systems – Home Assignment 1

Part 1 – SGD

1. Objective Function

$$\min_{x_u, y_i} \sum (r_{ui} - \hat{r}_{ui})^2 + \lambda \left[\sum_u \|x_u\|^2 + b_u^2 + \sum_i \|y_i\|^2 + b_i^2 \right]$$

2. Update Steps – Deriving the objective per parameter

$$\hat{r}_{ui} = x_u^T y_i + \mu + b_u + b_i$$

x_u

$$\frac{\partial}{\partial x_u} = -2 y_i \left[r_{ui} - \mu - b_u - b_i - x_u^T y_i \right] + 2 \lambda x_u \quad \text{v.2}$$

$$x_u^{\text{New}} = x_u^{\text{old}} + \eta (e_{ui} \cdot y_i - \lambda \cdot x_u)$$

y_i

$$\frac{\partial}{\partial y_i} = -2 x_u \left[r_{ui} - \mu - b_u - b_i - x_u^T y_i \right] + 2 \lambda y_i \quad \text{v.2}$$

$$y_i^{\text{New}} = y_i^{\text{old}} + \eta (e_{ui} \cdot x_u - \lambda \cdot y_i)$$

b_i

$$\frac{\partial}{\partial b_i} = -2 \cdot (r_{ui} - \mu - b_u - b_i - x_u^T y_i) + 2 \cdot \lambda \cdot b_i$$

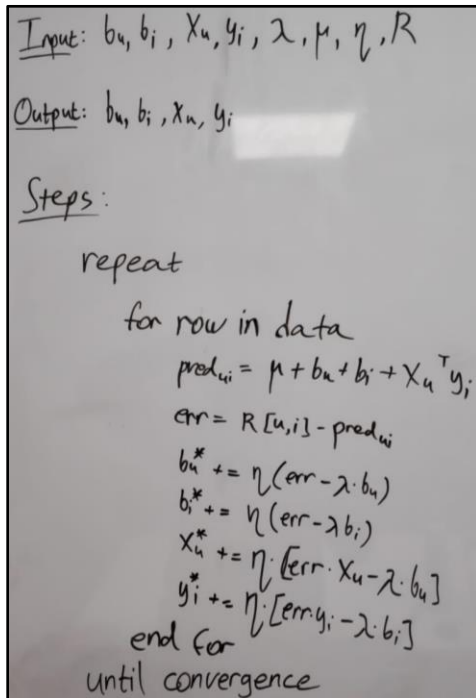
$$b_i^{\text{New}} = b_i^{\text{old}} + \eta \cdot (e_{ui} - \lambda \cdot b_i)$$

b_u

$$\frac{\partial}{\partial b_u} = -2 \cdot (r_{ui} - \mu - b_u - b_i - x_u^T y_i) + 2 \cdot \lambda \cdot b_u$$

$$b_u^{\text{New}} = b_u^{\text{old}} + \eta \cdot (e_{ui} - \lambda \cdot b_u)$$

3. Pseudo Code



Input: $b_u, b_i, X_u, y_i, \lambda, \mu, \eta, R$

Output: b_u, b_i, X_u, y_i

Steps:

```
repeat
  for row in data
     $\text{pred}_{ui} = \mu + b_u + b_i + X_u^T y_i$ 
     $\text{err} = R[u, i] - \text{pred}_{ui}$ 
     $b_u^* += \eta (\text{err} - \lambda \cdot b_u)$ 
     $b_i^* += \eta (\text{err} - \lambda \cdot b_i)$ 
     $X_u^* += \eta [\text{err} \cdot X_u - \lambda \cdot b_u]$ 
     $y_i^* += \eta [\text{err} \cdot y_i - \lambda \cdot b_i]$ 
  end for
until convergence
```

4. Hyper-parameters Tuning

For SGD we will tune the parameters:

- K – Dimensions
- Eta – Learning rate
- Lambda – Regularization rate (for each individual parameter)

5. Working with the validation set and convergence tests

After each iteration on the training set, we will check what is the RMSE on the validation set. Within the training set, we iterate over epochs. We defined convergence when the difference between epochs is smaller than 0.001 change or reaching the maximum number of epochs defined.

We use the validation set to determine which hyper-parameter is the most optimal to use in our final model.

6. Training the best\last model

By using the best model we found based on the training set and RMSE value on the validation set, we will then run it on both the training and validation sets combined. After convergence, we will use the model to run against the test set.

7. SGD Solution – main work items

In order to achieve a better result we made a few steps in the code:

- Clipping – changing prediction results that are lower than one or higher than five.
- Hyperparameters Tuning – we used random search method which sample different parameters values in each iteration, in order to find the most suitable parameters for our training set.

8. Results (RMSE, MAE, R^2)

RMSE is: 0.889745080993

MAE is: 0.6941523622081919

R^2 is: 0.363350125067

Part 2 – ALS

1. Objective Function

$$\min_{X_u, y_i} \sum (r_{ui} - \hat{r}_{ui})^2 + \lambda \left[\sum_u \|X_u\|^2 + b_u^2 + \sum_i \|y_i\|^2 + b_i^2 \right]$$

Exactly the same objective function as we used in the SGD model.

2. Update Steps – Deriving the objective per parameter

Preliminary step

$$\sum_D (r_{ui} - \hat{r}_{ui})^2 = \sum_D \left(r_{ui} - \mu - \overbrace{b_u}^a - \overbrace{b_i}^b - X_u^T y_i \right)^2 = \{a^2 - 2ab + b^2\}$$

$$\sum_D \left[(r_{ui} - \mu - b_u - b_i)^2 - 2 \cdot (r_{ui} - \mu - b_u - b_i) \cdot (X_u^T y_i) + X_u^T y_i \cdot y_i^T X_u \right] + \lambda \cdot \left(\sum_u \|X_u\|^2 + b_u^2 + \sum_i \|y_i\|^2 + b_i^2 \right)$$

X_u

$$\frac{\partial}{\partial X_u} = \sum_{u,i,r_{ui} \in D} [-2(r_{ui} - \mu - b_u - b_i)y_i + 2 \cdot y_i \cdot y_i^T X_u] + 2\lambda X_u = 0$$

$$\sum_D [y_i \cdot y_i^T X_u] + \lambda X_u = \sum_D [r_{ui} - \mu - b_u - b_i] \cdot y_i$$

$$\left[\sum_D [y_i y_i^T] + \lambda I \right] X_u = \sum_D [r_{ui} - \mu - b_u - b_i] y_i \quad \left(\left[\sum_D [y_i y_i^T] + \lambda I \right]^{-1} \right)$$

$$X_u^{New} = \left(\sum_D [y_i y_i^T] + \lambda I \right)^{-1} \cdot \sum_D [r_{ui} - \mu - b_u - b_i] y_i$$

y_i

$$\frac{\partial}{\partial y_i} = \sum_{u,i,r_{ui} \in D} [-2(r_{ui} - \mu - b_u - b_i)X_u + 2 \cdot X_u \cdot X_u^T y_i] + 2\lambda y_i = 0$$

$$\sum_D [X_u \cdot X_u^T y_i] + \lambda y_i = \sum_D [r_{ui} - \mu - b_u - b_i] \cdot X_u$$

$$\left[\sum_D [X_u X_u^T] + \lambda I \right] y_i = \sum_D [r_{ui} - \mu - b_u - b_i] X_u \quad \left(\left[\sum_D [X_u X_u^T] + \lambda I \right]^{-1} \right)$$

$$y_i^{New} = \left(\sum_D [X_u X_u^T] + \lambda I \right)^{-1} \cdot \sum_D [r_{ui} - \mu - b_u - b_i] X_u$$

bu

$$\begin{aligned}\frac{\partial}{\partial b_u} &= \sum_{D_u} 2(r_{ui} - \mu - b_i - x_u^T y_i + 2b_u) + 2 \cdot \lambda_{b_u} \cdot b_u = 0 \\ \sum_{D_u} b_u + \lambda_{b_u} \cdot b_u &= \sum_{D_u} (r_{ui} - \mu - b_i - x_u^T y_i) \\ (|D_u| + \lambda_{b_u}) b_u &= \sum_{D_u} (r_{ui} - \mu - b_i - x_u^T y_i) \cdot (|D_u| + \lambda_{b_u})^{-1} \\ b_u^{New} &= (|D_u| + \lambda)^{-1} \cdot \sum_{D_u} (r_{ui} - \mu - b_i - x_u^T y_i)\end{aligned}$$

bi

$$\begin{aligned}\frac{\partial}{\partial b_i} &= \sum_{D_i} 2(r_{ui} - \mu - b_u - x_u^T y_i + 2b_i) + 2 \cdot \lambda_{b_i} \cdot b_i = 0 \\ \sum_{D_i} b_i + \lambda_{b_i} \cdot b_i &= \sum_{D_i} (r_{ui} - \mu - b_u - x_u^T y_i) \\ (|D_i| + \lambda_{b_i}) b_i &= \sum_{D_i} (r_{ui} - \mu - b_u - x_u^T y_i) \cdot (|D_i| + \lambda_{b_i})^{-1} \\ b_i^{New} &= (|D_i| + \lambda)^{-1} \cdot \sum_{D_i} (r_{ui} - \mu - b_u - x_u^T y_i)\end{aligned}$$

3. Pseudo Code

Input: $b_u, b_i, x_u, y_i, \lambda, \mu, R$
Output: b_u, b_i, x_u, y_i
Steps:
 repeat
 for user in users: # fix i and estimate u
 $x_u^* = (\sum_D [y_i \cdot y_i^T] + \lambda \cdot I) \cdot (\sum_D r_{ui} - \mu - b_u - b_i) y_i$
 $b_u^* = \sum_{A_u} (r_{ui} - \mu - b_i - x_u^T y_i) \cdot (\|D_u\| + \lambda)^{-1}$
 for item in items: # fix u and estimate i
 $y_i^* = (\sum_D [x_u x_u^T] + \lambda \cdot I) \cdot (\sum_D r_{ui} - \mu - b_u - b_i) x_u$
 $b_i^* = \sum_{O_i} (r_{ui} - \mu - b_u - x_u^T y_i) \cdot (\|D_i\| + \lambda)^{-1}$
 and for s
 until convergence

4. Hyperparameters Tuning:

- K – Dimensions
- Lambda – Regularization rate.

5. Working with the validation set and convergence tests

Same as in SGD, after each iteration on the training set, we will check what is the RMSE on the validation set. We defined convergence when the difference between iteration is smaller than 0.001 change or reaching the maximum number of iterations defined.

6. Training the best\last model

Same as SGD. By using the best model, we found based on the training set, we will then run it on both the training and validation sets combined. After convergence, we will use the model to run against the test set.

7. ALS Solution – main work items:

- Clipping – changing prediction results that are lower than one or higher than five.
- Hyperparameters Tuning – we used random search method which sample different parameters values in each iteration, in order to find the most suitable parameters for our training set

8. Results (RMSE, MAE, R^2)

RMSE is: 0.9331830495897155

MAE is: 0.7141202638038082

R^2 is: 0.2996693236087469

9. Compare the ALS and SGD solutions in terms of implementation, training and quality.

SGD

Implementation – We first initialized the data into a pandas dataframe. Afterwards we created three main data structures to be used in the rest of the stages – two pandas series for the users and items, which was used to index all the available entities. In addition, we used pandas pivot table to recreate the original data into a matrix of users, items and their corresponding ratings.

Training – In each iteration we updated the latent matrices and the model parameters for all the users and items. We stopped iterating if the convergence condition has been satisfied, as described in the SGD section.

Quality – We will compare the models based on running time and efficiency. The SGD running times until the model converged were slightly longer, however the result received at the end were slightly better than the ones we received using the ALS model.

ALS

Implementation – We used the same structure as we did in the SGD model.

Training – On each iteration we first fixed the items latent matrix and updated the users parameters according to the fixed values, as a second step we did the exact same thing just changing the users to be fixed instead of the items.

Quality – As described in the SGD quality section above, the ALS model converged faster than the SGD one, but yielded a bit less results in terms of performance.