

# Assignment 3

Authors:

Name	ID
Sangi Manish Rao	2015A7PS0235H
Tanuj Gupta	2015A7PS0159H
Krishna Bharadwaj	2015A7PS0076H
Ashrith Grandi	2015A7PS0285H

Code Files:

1. cur.py
2. cur\_energy.py
3. svd.py
4. svd\_energy.py
5. collab\_filtering.py
6. collab\_filtering\_baseline.py

Working:

Described in README file

Language:

Python

Packages Used:

1. Timeit
2. Numpy
3. Math

Project Description:

In this project we have implemented three major recommender systems which are Collaborative Systems, SVD and CUR along with some modifications.

## Collaborative Systems:

Pearson correlation coefficient was used for similarity measure.

Formula:

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

User-User collaborative filtering is used. The following formula is used to predict a user's movie rating.

$$r_{xi} = \frac{\sum_{j \in N(i; x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i; x)} s_{ij}}$$

$s_{ij}$ ... similarity of items  $i$  and  $j$   
 $r_{xj}$ ... rating of user  $u$  on item  $j$   
 $N(i; x)$ ... set items rated by  $x$  similar to  $i$

## Baseline estimate Collaborative-filtering:

Pearson correlation coefficient was used for similarity measure.

User-User collaborative filtering is used. The following formula is used to predict a user's movie rating

- Estimate rating  $r_{xi}$  as the weighted average:

$$r_{xi} = b_{xi} + \frac{\sum_{j \in N(i;x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i;x)} s_{ij}}$$

baseline estimate for  $r_{xi}$

$$b_{xi} = \mu + b_x + b_i$$

- $\mu$  = overall mean movie rating
- $b_x$  = rating deviation of user  $x$   
= (avg. rating of user  $x$ ) -  $\mu$
- $b_i$  = rating deviation of movie  $i$

## Singular Value decomposition:

Given matrix is decomposed to these components.

$$\mathbf{A}_{[m \times n]} = \mathbf{U}_{[m \times r]} \Sigma_{[r \times r]} (\mathbf{V}_{[n \times r]})^T$$

**A: Input data matrix**

–  $m \times n$  matrix (e.g.,  $m$  users,  $n$  movies)

**U: Left singular vectors**

–  $m \times r$  matrix ( $m$  users,  $r$  concepts)

**$\Sigma$ : Singular values**

–  $r \times r$  diagonal matrix (strength of each 'concept')  
( $r$ : rank of the matrix  $\mathbf{A}$ )

**V: Right singular vectors**

–  $n \times r$  matrix ( $n$  movies,  $r$  concepts)

### SVD with 90% energy:

To find out this, the squares of the diagonal elements of Sigma are summed, until the sum doesn't exceed 90% of the total energy (which is sum of squares of all the diagonal elements). At the point it breaks, the index is noted and dimensionality reduction is done accordingly.

# CUR:

## ■ Sampling columns (similarly for rows):

Total length of all the columns

**Input:** matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , sample size  $c$

**Output:**  $\mathbf{C}_d \in \mathbb{R}^{m \times c}$

1. for  $x = 1 : n$  [column distribution]
2.  $P(x) = \sum_i \mathbf{A}(i, x)^2 / \sum_{i,j} \mathbf{A}(i, j)^2$
3. for  $i = 1 : c$  [sample columns]
4. Pick  $j \in 1 : n$  based on distribution  $P(j)$
5. Compute  $\mathbf{C}_d(:, i) = \mathbf{A}(:, j) / \sqrt{cP(j)}$

Note this is a randomized algorithm, same column can be sampled more than once

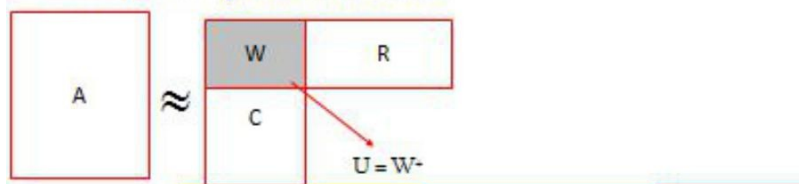
- Let  $\mathbf{W}$  be the “intersection” of sampled columns  $\mathbf{C}$  and rows  $\mathbf{R}$

- Let SVD of  $\mathbf{W} = \mathbf{X} \mathbf{Z} \mathbf{Y}^T$

- Then:  $\mathbf{U} = \mathbf{W}^+ = \mathbf{Y} \mathbf{Z}^+ \mathbf{X}^T$

- $\mathbf{Z}^+$ : reciprocals of non-zero singular values:  $Z^+_{ii} = 1 / Z_{ii}$

- $\mathbf{W}^+$  is the “pseudoinverse”



dense but small

$$\text{CUR: } \mathbf{A} = \mathbf{C} \mathbf{U} \mathbf{R}$$

Huge but sparse

Big but sparse

CUR with 90% energy:

To compute this  $\mathbf{U} = \text{pseudoinverse}(\mathbf{W})$

Where W is the SVD with 90% energy of the matrix which is intersection of C,R.

## Error-calculations:

### Root mean square error (rmse):

▪ Root-mean-square error (RMSE)

▪  $\sqrt{\sum_{\substack{xi \\ \text{of } x \text{ on } i}} (r_{xi} - r_{xi}^*)^2}$  where  $r_{xi}$  is predicted,  $r_{xi}^*$  is the true rating

### Precision on top k:

k=50

Relevance = 3

If the number of user ratings considered are less than k

Value for each user = (Number of predicted user ratings greater than relevance)/(Number of user ratings actually greater than relevance)

Precision on top k = Mean of the values for all relevance

### Rank Correlation:

▪ Spearman's *correlation* between system's and user's complete rankings

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

where  $d_i$  = absolute difference in user and predicted ratings  
 $n$  = number of ratings

### Comparisons:

Recommender System	RMSE	Precision on Top k(k=3, rel=2)	Spearman Rank Correlation	Time taken
Collaborative	1.31	0.0010604453870625664	0.99	66.46347651415039
Collaborative with Baseline Approach	1.024	0.003534817956875221	0.99	68.66220501824604
SVD	0.423	0.3716861081654297	0.99	236.37811649303268
SVD with 90% Energy	0.24	0.3891834570519618	0.99	236.53716842969106
CUR	Depends on c,r but average is around 100	0.1378579003181332	0.98	15.053157166368827
CUR with 90% energy	Depends on c,r but average is around 100	0.10763520678685047	0.98	11.98424341847331

## Data Set:

Small: 100,000 ratings applications applied to 1046 movies by 943 users. Last updated 10/2016.