

# INFORMATION RETRIEVAL

CS F469

## RECOMMENDER SYSTEMS

Pranjal Gupta

2017A7PS0124H

Daksh Yashlaha

2017A7PS0218H

Satyam Mani

2017A7PS0277H

Ujjwal Raizada

2017A7PS1398H

## Few important details

- Dataset used: MovieLens dataset
- Strict and easy raters are handled properly using pearson correlation
- CUR details
  - Rows -
  - Columns -

<b>Recommender System Technique</b>	<b>Root Mean Squared Error (RMSE)</b>	<b>Mean Average Error (MAE)</b>	<b>Time taken for prediction</b>
Collaborative	1.158105984859306	0.59734323668047	12.69501471519s
Collaborative along with Baseline approach	0.59405799171546	0.511784030852735	8.78501643517s
SVD	0.89123464828352	0.854693215875624	12.8539641275s
SVD with 90 per cent energy	0.93576824631325	0.902395864782184	10.6939641351s
CUR	0.98564682842824	0.945845754924922	12.641667522s
CUR with 90 per cent energy	1.04542854655445	0.965846427454562	12.2184527496s
Latent factor model	0.89146426426542	0.862584645656546	12.8454651465s

- Language Used : python3
- Libraries used
  - Numpy
  - Pandas
  - Scipy
  - Sklearn

## Collaborative Filtering and filtering with baseline approach

The process of identifying similar users and recommending what similar users like is called collaborative filtering.

Similarity Measure used: Pearson correlation

Major Data-Structures used:

- df : pandas dataframe
  - index (along the row): user
  - label (columns) : movies
- utility\_matrix : 2d numpy array
  - row - user\_id
  - Column - movie\_id
- Users\_map: { key, value }
  - key : user\_no
  - value : iterator in the range(len(user\_map))
- Movie\_map: { key, value }
  - key : movie\_no
  - value : iterator in the range(len(movie\_map))

\*\*\* Same data-structures were followed for the baseline approach \*\*\*

## Singular Value Decomposition

**SVD:** SVD is calculated for any given  $M \times N$  matrix by decomposing that matrix into three component matrices defined as follows:

- Matrix  $U$ : An  $M \times r$  column orthonormal matrix where  $r$  is the rank of the original matrix.
- Matrix  $\Sigma$ : An  $r \times r$  diagonal matrix containing the singular values of the original matrix.
- Matrix  $V$ : An  $r \times N$  column orthonormal matrix where  $r$  is the rank of the original matrix.

The original matrix can then be approximated by calculating the product of these matrices as follows:

$$\text{Approximation of original matrix} = U * \Sigma * V^T$$

### Dimensionality Reduction:

If the rank  $r$  of the matrix is substantially smaller than  $N$ , we can intuitively see that the combined sizes of the component matrices will be smaller than that of the original matrix, thereby providing a valuable reduction in dimensionality for computational purposes.

**90% energy rule:** Retain enough singular values to make up 90% of the energy. That is, the sum of squares of the retained singular values should be at least 90% of the sum of squares of all the singular values.

## CUR Method and CUR with 80% - 90% energy

Utility matrix is divided into three matrices C, U, R

Major Data Structure:

- Matrix C: An  $M \times r$  column orthonormal matrix where  $r$  is the rank of the original matrix.
- Matrix U: An  $r \times r$  diagonal matrix containing the singular values of the original matrix.
- Matrix R: An  $r \times N$  column orthonormal matrix where  $r$  is the rank of the original matrix.

The original matrix can then be approximated by calculating the product of these matrices as follows:

Approximation of original matrix =  $C * U * R^T$

## LATENT FACTOR MODEL

In latent factor based method, we only feed user's history and we need not to define descriptors or factors. The algorithm will find the hidden factors that influence the user's preference like brand of product, price of product and so on for us. Here, no product description is required, only user's history is good.

Major Data Structures Used:

- utility\_matrix: 2d array
  - row - user\_id
  - columns - movie\_id
- user\_factor\_matrix : 2d array
  - Computed after decomposition of utility matrix into latent factors
- product\_factor\_matrix : 2d array
  - Second result due to decomposition of utility matrix into latent factors