# Clustering Based Recommendation Systems

Capstone Project: Team #14
1. Rahul
2. Siddhant
3. Himangshu
4. Sushma
5. Krunal
6. Harsha

# Introduction

# INTRODUCTION



- Recommendations systems give two kinds of feedbacks:
  - Explicit feedback from Users
  - Implicit feedback from associated with user data
- Associating all the users and products for generating recommendations are often not meaningful

# Dataset

# Data Set



- MovieLens Dataset
- 100,000 Ratings (1 - 5)
- Available in different granularities.
- Structure:
  - 943 Users on 1682 Movies
  - Each user has rated at least 20 movies
  - Simple  demographics info of users.

# BI Use Case

# BI use case : Problems with Original Algorithms

- The raw algorithms run on the whole database of users and items.

- Two typical problems:

    - Exponentially huge size of User-Item matrix.

    - Doesn't leverage the association, similarity and dissimilarity

        between users and items.

# Approach

# Clustering- The Concept

- Significance:
  - Speed
  - Size
  - Meaning of Recommendation

- Types:
  - User Clustering
  - Item Clustering

# Optimist, Pessimist and Neutral Users

- Normally done using Demographics, but not good.

- Using user prior recommendation information.

- Set α and β, as defined in the paper.

- Select UL, UH and UN

- Find the leaders.

# Similarity Matrices

- Baseline Matrices PCC and Cosine.

- New matrix UPS.

- $$sim(a, b)^{UPS} = \exp\left(-\frac{\sum_{i \in I_{ab}} |r_{ai} - r_{bi}|}{|I_{ab}|} \times |\bar{r}_a - \bar{r}_b|\right)$$
$$\times \frac{|I_a| \cap |I_b|}{|I_a| \cup |I_b|}, \tag{6}$$

# Clustering

$u \in C_o, \quad \text{if } u \text{ satisfies } sim(u, c_o) > sim(u, c_p), \text{ and}$
$\quad sim(u, c_o) > sim(u, c_n);$
$u \in C_p, \quad \text{if } u \text{ satisfies } sim(u, c_p) > sim(u, c_o), \text{ and}$
$\quad sim(u, c_p) > sim(u, c_n);$
$u \in C_n, \quad \text{if } u \text{ satisfies } sim(u, c_n) > sim(u, c_o), \text{ and}$
$\quad sim(u, c_n) > sim(u, c_p).$

$$c_o = u \Leftarrow \arg\max_u |I_u|, \tag{4}$$

where $\forall u \in U_h$, $I_u = \{i \in I | r_{ui} \neq ?\}$. Clustering center $c_p$ of pessimistic user group can be uniquely determined, as follow:

$$c_p = u \Leftarrow \arg\max_u |I_u|, \tag{5}$$

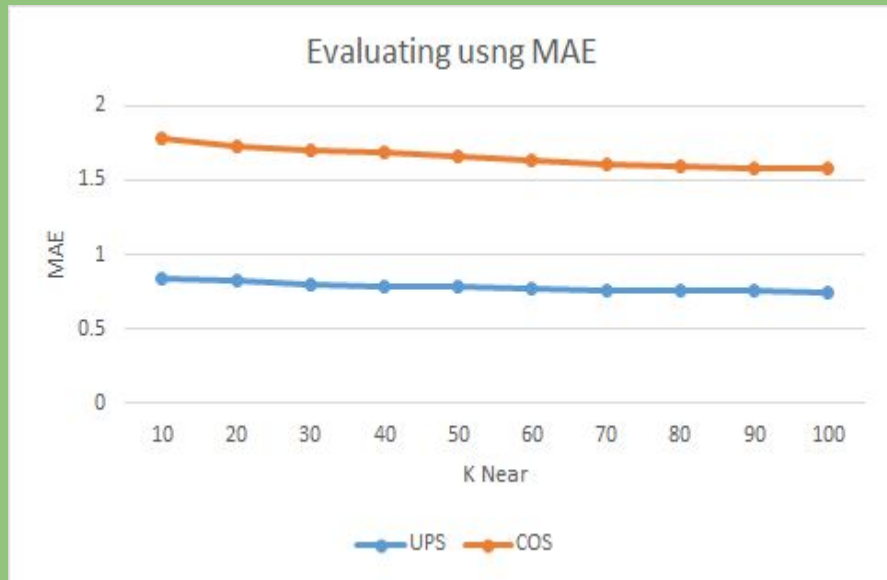where $\forall u \in U_l$, $I_u = \{i \in I | r_{ui} \neq ?\}$.

# Recommendations

**Definition 5** Suppose $U_{nei}$ is the neighbor set of active user $t$, $|U_{nei}| = k$, and $U_{nei} \subset U$. If $t \in U_o$, then $U_{nei} \subset U_o \cup U_n$. If $t \in U_p$, then $U_{nei} \subset U_p \cup U_n$. If $t \in U_n$, then $U_{nei} \subset U_n$.

$$p_{ti} = \bar{r}_t + \frac{\sum_{u \in U_{nei}} sim^{UPS}(t, u) \times (r_{ui} - \bar{r}_u)}{\sum_{u \in U_{nei}} |sim^{UPS}(t, u)|},$$

# Evaluation and Observation

# Evaluation and observations



Evaluating usng MAE

- $$MAE = \frac{1}{\#U} \sum_{u \in U} \frac{\sum_{i \in I_u} |p_{ui} - r_{ui}|}{\#I_u}$$

- UPS similarity metrics as compared to Cosine similarity metrics

- Optimal value of K

# Conclusion and Future Work

# Conclusion and Future Work

- Introducing the Bias into the recommendation system.

- How will it behave if done for items?

- Refine the criteria of optimism and pessimism.

- The people who give bad ratings on very famous and enjoyed items are the pessimists!