

Item Based KNN Collaborative Filtering

J. Milmore

University of Massachusetts Amherst
Department of Mathematics and Statistics

April 6, 2020

Item Based CF

Memory Based Collaborative Filtering: Predict a rating for user u and movie m by using previous ratings

- ▶ Item Based: $\text{sim}(m, m')$

Problem Formulation

Rating Matrix

- ▶ Rows represent ratings given to each movie
- ▶ Columns represent ratings given by each user
- ▶ $R_{m,u}$ represents rating given to movie m by user u

$$R = \begin{matrix} & u_1 & u_2 & u_* \\ \begin{matrix} m_* \\ m_2 \\ m_3 \\ m_4 \\ m_5 \end{matrix} & \begin{pmatrix} 2 & 5 & ? \\ 3 & ? & 3 \\ ? & 4 & 2 \\ 3 & ? & 4 \\ 2 & 5 & ? \end{pmatrix} \end{matrix}$$

Goal

- ▶ Predict R_{m_*,u_*} , the missing rating for user u_* and movie m_* using the existing ratings

Method

- Filter rating matrix to just contain rows of movies rated by u_* , along with the movie of interest, m_*

$$\begin{array}{c} \begin{array}{ccc} & u_1 & u_2 & u_* \\ m_* & \left(\begin{array}{ccc} 2 & 5 & ? \\ 3 & ? & 3 \\ ? & 4 & 2 \\ 3 & ? & 4 \\ 2 & 5 & ? \end{array} \right) \end{array} \end{array} \rightarrow \begin{array}{c} \begin{array}{ccc} & u_1 & u_2 & u_* \\ m_* & \left(\begin{array}{ccc} 2 & 5 & ? \\ 3 & ? & 3 \\ ? & 4 & 2 \\ 3 & ? & 4 \end{array} \right) \end{array} \end{array}$$

Method (Cont.)

- Replace missing ratings by the average rating of the user

$$\begin{array}{c} u_1 \quad u_2 \quad u_* \\ m_* \begin{pmatrix} 2 & 5 & ? \end{pmatrix} \\ m_2 \begin{pmatrix} 3 & ? & 3 \end{pmatrix} \\ m_3 \begin{pmatrix} ? & 4 & 2 \end{pmatrix} \\ m_4 \begin{pmatrix} 3 & ? & 4 \end{pmatrix} \end{array} \rightarrow \begin{array}{c} u_1 \quad u_2 \quad u_* \\ m_* \begin{pmatrix} 2 & 5 & \bar{u}_* \end{pmatrix} \\ m_2 \begin{pmatrix} 3 & \bar{u}_2 & 3 \end{pmatrix} \\ m_3 \begin{pmatrix} \bar{u}_1 & 4 & 2 \end{pmatrix} \\ m_4 \begin{pmatrix} 3 & \bar{u}_2 & 4 \end{pmatrix} \end{array}$$

Method (Cont.)

- ▶ Adjust each rating by the average rating of the user

$$\begin{array}{c} u_1 \quad u_2 \quad u_* \\ m_* \begin{pmatrix} 2 & 5 & \bar{u}_* \end{pmatrix} \\ m_2 \begin{pmatrix} 3 & \bar{u}_2 & 3 \end{pmatrix} \\ m_3 \begin{pmatrix} \bar{u}_1 & 4 & 2 \end{pmatrix} \\ m_4 \begin{pmatrix} 3 & \bar{u}_2 & 4 \end{pmatrix} \end{array} \rightarrow \begin{array}{c} u_1 \quad u_2 \quad u_* \\ m_* \begin{pmatrix} 2 - \bar{u}_1 & 5 - \bar{u}_2 & 0 \end{pmatrix} \\ m_2 \begin{pmatrix} 3 - \bar{u}_1 & 0 & 3 - \bar{u}_* \end{pmatrix} \\ m_3 \begin{pmatrix} 0 & 4 - \bar{u}_2 & 2 - \bar{u}_* \end{pmatrix} \\ m_4 \begin{pmatrix} 3 - \bar{u}_1 & 0 & 4 - \bar{u}_* \end{pmatrix} \end{array}$$

We will further refer to this *adjusted* rating matrix simply as R

- ▶ Find the k movies most **similar** to m_* and **predict** (m_*, u_*) to be the weighted average of the ratings given by u_* to these k movies

Scoring Similar Movies

Adjusted Cosine Similarity Metric

$$\text{sim}(m_1, m_2) = \frac{\sum_{u \in A} (R_{m_1, u} - \bar{R}_u)(R_{m_2, u} - \bar{R}_u)}{\sum_{u \in A} (R_{m_1, u} - \bar{R}_u)^2 \sum_{u \in A} (R_{m_2, u} - \bar{R}_u)^2}$$

Where $A = U(m_1) \cap U(m_2) = \{u: u \text{ has rated both } m_1 \text{ and } m_2 \}$

Calculating Similarity

- ▶ Want to calculate the similarity between m_* and all

$$m_i \in M(u_*)$$

- ▶ Break into subproblems. For each $m_i \in M(u_*)$

1. Calculate $\sum_{u \in A} (R_{m_*,u} - \bar{R}_u)(R_{m_i,u} - \bar{R}_u)$

2. Calculate $\sum_{u \in A} (R_{m_*,u} - \bar{R}_u)^2$

3. Calculate $\sum_{u \in A} (R_{m_i,u} - \bar{R}_u)^2$

Step 1

- ▶ Calculate $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)(R_{m_i, u} - \bar{R}_u)$
- ▶ Multiply adjusted rating matrix by row vector of ratings for m_*

$$\begin{array}{c} m_* \\ m_1 \\ m_2 \\ m_3 \end{array} \begin{array}{ccc} u_1 & u_2 & u_* \\ \left(\begin{array}{ccc} 2 - \bar{u}_1 & 5 - \bar{u}_2 & 0 \\ 3 - \bar{u}_1 & 0 & 3 - \bar{u}_* \\ 0 & 4 - \bar{u}_2 & 2 - \bar{u}_* \\ 3 - \bar{u}_1 & 0 & 4 - \bar{u}_* \end{array} \right) \end{array} \times (2 - \bar{u}_1 \quad 5 - \bar{u}_2 \quad 0)^T$$

Step 1 (Cont.)

Resulting Vector

$$\begin{matrix} m_*, m_* \\ m_1, m_* \\ m_2, m_* \\ m_3, m_* \end{matrix} \begin{pmatrix} (2 - \bar{u}_1)^2 + (5 - \bar{u}_2)^2 \\ (3 - \bar{u}_1)(2 - \bar{u}_1) \\ (4 - \bar{u}_2)(5 - \bar{u}_2) \\ (3 - \bar{u}_1)(2 - \bar{u}_1) \end{pmatrix}$$

- ▶ Each entry is the sum of products of adjusted ratings between m_* and $m_i \in M(u_*)$
- ▶ The only movie ratings considered are from users that have rated both m_* and m_i
- ▶ Each entry equal to $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)(R_{m_i, u} - \bar{R}_u)$

Step 2

Calculate $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)^2$

- ▶ Transpose the adjusted rating matrix: R^T
- ▶ Transform the vector of ratings m_* into a matrix of same shape as R^T such that each column represents m_* : M_*
- ▶ Multiply M_* by R^T elementwise, then divide by R^T elementwise

Step 2 (Cont.)

$$M_* \begin{array}{c} u_1 \\ u_2 \\ u_* \end{array} \begin{array}{cccc} m_* & m_* & m_* & m_* \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 \\ 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 \\ 0 & 0 & 0 & 0 \end{array} \right) * \end{array}$$

$$R^T \begin{array}{c} u_1 \\ u_2 \\ u_* \end{array} \begin{array}{cccc} m_* & m_1 & m_2 & m_3 \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 3 - \bar{u}_1 & 0 & 3 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 4 - \bar{u}_2 & 0 \\ 0 & 3 - \bar{u}_* & 2 - \bar{u}_* & 4 - \bar{u}_* \end{array} \right) \div \end{array}$$

$$R^T \begin{array}{c} u_1 \\ u_2 \\ u_3 \end{array} \begin{array}{cccc} m_* & m_1 & m_2 & m_3 \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 3 - \bar{u}_1 & 0 & 3 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 4 - \bar{u}_2 & 0 \\ 0 & 3 - \bar{u}_* & 2 - \bar{u}_* & 4 - \bar{u}_* \end{array} \right) \end{array}$$

Step 2 (Cont.)

$$= \begin{matrix} & m_* \sim m_* & m_* \sim m_1 & m_* \sim m_2 & m_* \sim m_3 \\ \begin{matrix} u_1 \\ u_2 \\ u_* \end{matrix} & \begin{pmatrix} 2 - \bar{u}_1 & 2 - \bar{u}_1 & 0 & 2 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 5 - \bar{u}_2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

In the resulting matrix, entry (u_i, m_i) is the adjusted rating of user u_i on movie m_* such that u_i has rated both m_i and m_* , 0 otherwise.

- Now multiply the row vector of ratings for movie m_* by the resulting matrix

Step 2 (Cont.)

$$(2 - \bar{u}_1 \quad 5 - \bar{u}_2 \quad 0)^* \begin{matrix} & m_* \sim m_* & m_* \sim m_1 & m_* \sim m_2 & m_* \sim m_3 \\ \begin{matrix} u_1 \\ u_2 \\ u_* \end{matrix} & \begin{pmatrix} 2 - \bar{u}_1 & 2 - \bar{u}_1 & 0 & 2 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 5 - \bar{u}_2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

$$= ((2 - \bar{u}_1)^2 + (5 - \bar{u}_2)^2 \quad (2 - \bar{u}_1)^2 \quad (5 - \bar{u}_2)^2 \quad (2 - \bar{u}_1)^2)$$

Resulting Vector

- ▶ Each entry represents the sum of squares of the adjusted ratings of m_* given by users that have rated both m_* and m_i
- ▶ Each entry equal to $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)^2$

Step 3

Calculate $\sum_{u \in A} (R_{m_i, u} - \bar{R}_u)^2$

- ▶ Transform the vector of ratings m_* into a matrix of same shape as R^T such that each column represents m_* : M_*
- ▶ Transpose the adjusted rating matrix: R^T
- ▶ Multiply R^T by M_* elementwise, then divide by M_* elementwise

Step 3 (Cont.)

$$R_T \begin{array}{c} u_1 \\ u_2 \\ u_* \end{array} \begin{array}{cccc} m_* & m_1 & m_2 & m_3 \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 3 - \bar{u}_1 & 0 & 3 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 4 - \bar{u}_2 & 0 \\ 0 & 3 - \bar{u}_* & 2 - \bar{u}_* & 4 - \bar{u}_* \end{array} \right) * \end{array}$$

$$M_* \begin{array}{c} u_1 \\ u_2 \\ u_* \end{array} \begin{array}{cccc} m_* & m_* & m_* & m_* \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 \\ 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 \\ 0 & 0 & 0 & 0 \end{array} \right) \div \end{array}$$

$$M_* \begin{array}{c} u_1 \\ u_2 \\ u_* \end{array} \begin{array}{cccc} m_* & m_* & m_* & m_* \\ \left(\begin{array}{cccc} 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 & 2 - \bar{u}_1 \\ 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 & 5 - \bar{u}_2 \\ 0 & 0 & 0 & 0 \end{array} \right) * \end{array}$$

Step 3 (Cont.)

$$= \begin{matrix} & m_* \sim m_* & m_1 \sim m_* & m_2 \sim m_* & m_3 \sim m_* \\ \begin{matrix} u_1 \\ u_2 \\ u_* \end{matrix} & \begin{pmatrix} 2 - \bar{u}_1 & 3 - \bar{u}_1 & 0 & 3 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 4 - \bar{u}_2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{matrix}$$

In the resulting matrix, entry (u_i, m_i) is the adjusted rating of user u_i on movie m_i such that u_i has rated both m_i and m_* , 0 otherwise.

- Now multiply the rating matrix by the above matrix

Step 3 (Cont.)

$$\begin{array}{c} u_1 \quad u_2 \quad u_* \\ m_* \\ m_1 \\ m_2 \\ m_3 \end{array} \begin{pmatrix} 2 - \bar{u}_1 & 5 - \bar{u}_2 & 0 \\ 3 - \bar{u}_1 & 0 & 3 - \bar{u}_* \\ 0 & 4 - \bar{u}_2 & 2 - \bar{u}_* \\ 3 - \bar{u}_1 & 0 & 4 - \bar{u}_* \end{pmatrix} \times$$

$$\begin{array}{c} m_* \sim m_* \quad m_1 \sim m_* \quad m_2 \sim m_* \quad m_3 \sim m_* \\ u_1 \\ u_2 \\ u_* \end{array} \begin{pmatrix} 2 - \bar{u}_1 & 3 - \bar{u}_1 & 0 & 3 - \bar{u}_1 \\ 5 - \bar{u}_2 & 0 & 4 - \bar{u}_2 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Step 3 (Cont.)

$$\begin{array}{c} m_* \\ m_1 \\ m_2 \\ m_3 \end{array} \begin{pmatrix} m_* & m_1 & m_2 & m_3 \\ (2 - \bar{u}_1)^2 + (5 - \bar{u}_2)^2 & - & - & - \\ - & (3 - \bar{u}_1)^2 & - & - \\ - & - & (4 - \bar{u}_2)^2 & - \\ - & - & - & (3 - \bar{u}_1)^2 \end{pmatrix}$$

Resulting Matrix

- ▶ Entries along the diagonal represent the sum of squares of the adjusted ratings of m_i given by users that have rated both m_* and m_i
- ▶ Entries along the diagonal represent $\sum_{u \in A} (R_{m_i, u} - \bar{R}_u)^2$

Prediction

- ▶ Able to calculate $\text{sim}(m_*, m_i)$ for all $m_i \in M(u_*)$ in a few simple matrix operations
- ▶ For all $m_i \in M(u_*)$
 - ▶ Step 1 provides us $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)(R_{m_i, u} - \bar{R}_u)$
 - ▶ Step 2 provides us $\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)^2$
 - ▶ Step 3 provides us $\sum_{u \in A} (R_{m_i, u} - \bar{R}_u)^2$ along the diagonals of a matrix
 - ▶ Combine to get $\text{sim}(m_*, m_i) = \frac{\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)(R_{m_i, u} - \bar{R}_u)}{\sum_{u \in A} (R_{m_*, u} - \bar{R}_u)^2 \sum_{u \in A} (R_{m_i, u} - \bar{R}_u)^2}$
- ▶ Now can choose the k movies most similar to m_* (that are rated by u_*) to predict R_{m_*, u_*}

Prediction

Predict rating of user u for movie m by taking weighted average of the ratings given by the k movies most similar to m that have rating from u .

$$P_{m,u} = \frac{\sum_{m' \in N_u^K(m)} R_{m',u} \text{sim}(m, m')}{\sum_{m' \in N_u^K(m)} |\text{sim}(m, m')|}$$

$N_u^K(m) = \{m' : m' \text{ belongs to the } k \text{ most similar movies of } m \text{ and } u \text{ has rated } m'\}$