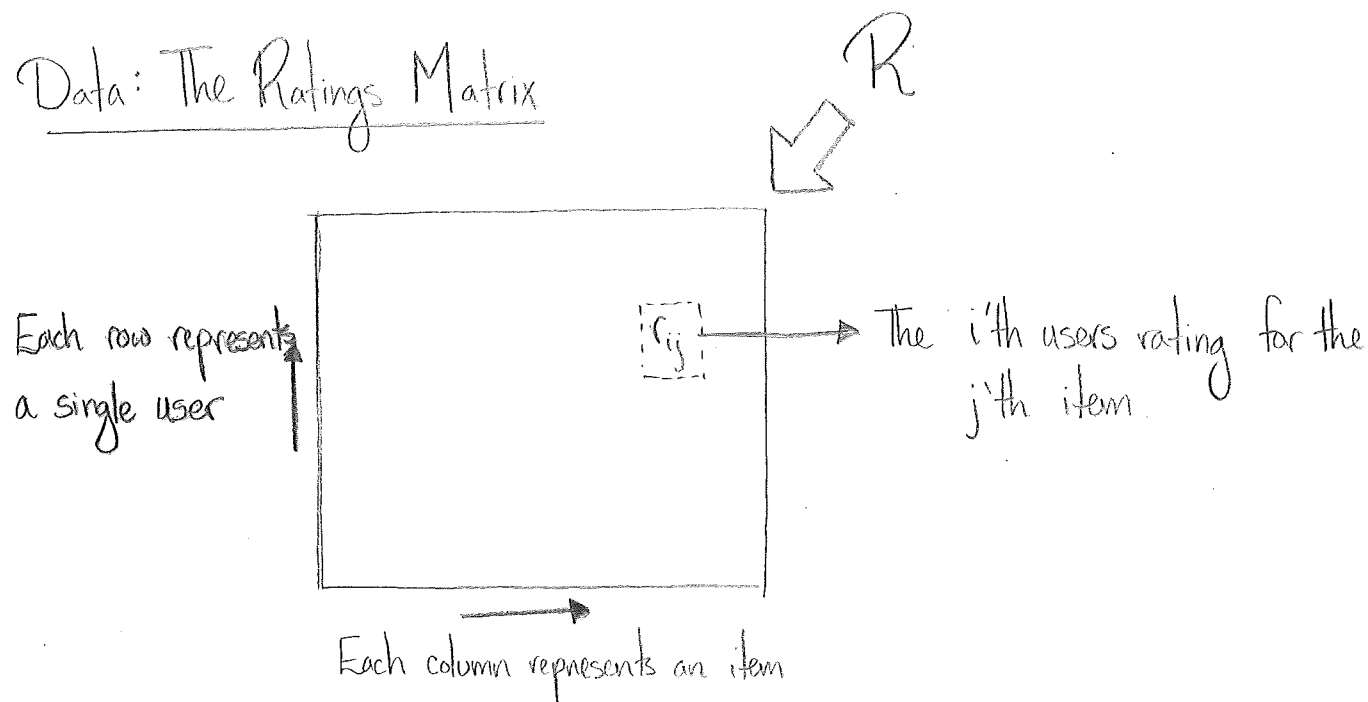


# Recommendation Systems

Morning: Similarity based / Collaborative Filtering

Afternoon: Matrix Factorization

## Data: The Ratings Matrix



Usually, most entries are missing, since most users have not rated most items!

Ratings can be explicit or implicit

Explicit: User supplies ratings for items

Implicit: User consumes selected items,

measure of consumption is taken as a rating.

Goal: Predict missing ratings!

$r_{ij}$ : Actual rating of user  $i$  for item  $j$

$\hat{r}_{ij}$ : Predicted rating of user  $i$  for item  $j$ .

# Similarity Based

## User - User Hypothesis :

Similar users tend to give similar ratings to a single product.

## Item - Item Hypothesis :

A single user will tend to give similar ratings to similar products.

Need : A way to measure similarity between users and/or items.

Then : We can predict ratings as a weighted average of the actual ratings

- of similar users (for a fixed item)
- of similar products (for a fixed user)

## User - User :

$$\hat{r}_{ij} = \frac{\sum_{\substack{\text{users } u \text{ that have} \\ \text{rated item } j}} \overbrace{\text{user-user similarity}}^{\text{sim}(i, u)} r_{uj}}{\sum_u \text{sim}(i, u)}$$

↑ rating for fixed item j.

## Discussions :

- ① What are the possible benefits / drawbacks of both approaches.
- ② How may we make these calculations more efficient?

## Item - Item :

$$\hat{r}_{ij} = \frac{\sum_{\substack{\text{items } t \text{ that are} \\ \text{rated by user } i}} \overbrace{\text{item-item similarity}}^{\text{sim}(j, t)} r_{it}}{\sum_t \text{sim}(j, t)}$$

## Similarity Measures :

### Requirements :

- ①  $\text{sim}(a, b)$  is between (inclusive) zero and one.
- ②  $\text{sim}(a, b) = 0$  means "a and b are not at all similar".
- ③  $\text{sim}(a, b) = 1$  means "a and b are as similar as possible".

Similarity between users/items is based off the rows (user-user) or columns (item-item) of the rating matrix.

### Cosine Similarity :

$$\cos(\theta_{\vec{a}, \vec{b}}) = \frac{\vec{a}}{|\vec{a}|} \cdot \frac{\vec{b}}{|\vec{b}|}$$

Notation

$\vec{a}, \vec{b}$  : rows (user-user) or columns (item-item) of the rating matrix  $R$ .

$$\text{cos-sim}(\vec{a}, \vec{b}) = \frac{1}{2} + \frac{1}{2} \cos(\theta_{\vec{a}, \vec{b}})$$

### Pearson - Correlation Similarity

$$\text{corr}(\vec{a}, \vec{b}) = \frac{\text{cov}(\vec{a}, \vec{b})}{\text{sd}(\vec{a}) \text{sd}(\vec{b})}$$

$$\text{corr-sim}(\vec{a}, \vec{b}) = \frac{1}{2} + \frac{1}{2} \text{corr}(\vec{a}, \vec{b})$$

### Jaccard - Similarity

$R$  must be binary : user consumed item or not?

$$\text{jacc-sim}(\vec{a}, \vec{b}) = \frac{\# \text{ of items consumed by both } a \text{ and } b}{\# \text{ of items consumed by either } a \text{ or } b}$$

### Discussion :

When are these large ( $\approx 1$ )

When are these small ( $\approx 0$ )

What should we do with missing ratings when computing similarity?