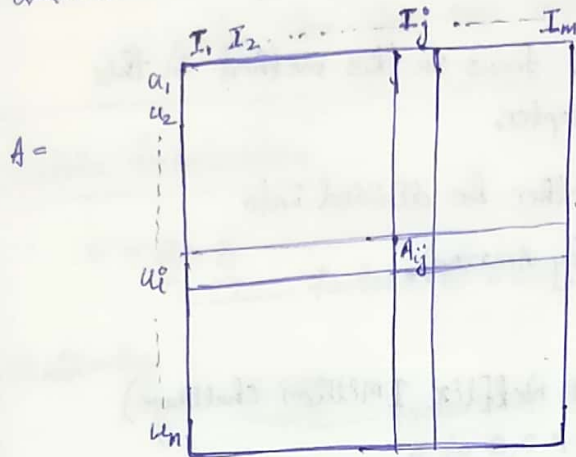


# Recommendation System

## Formulation

Dataset = A



$n$ : users  
 $m$ : items.

$A_{ij}$ : rating given by  $u_i$  on  $I_j$  (1-5)  
or  
 $u_i$  watched  $I_j$  (1 or 0).

\* very sparse matrix as any user would have rated around 5-10 items.

## Task of RS:-

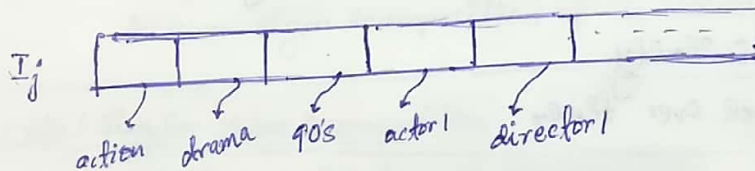
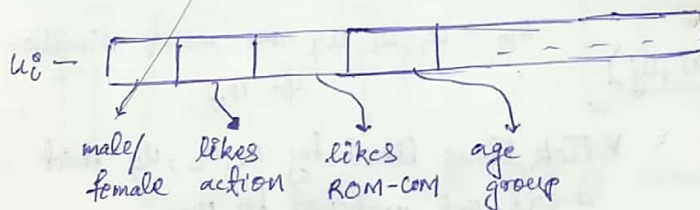
$u_i$ :  $I_1, I_2, I_7, I_8$  — items rated by user  $u_i$ .  
↳ recommendation of new items based on these ratings??

RS algorithms can be broadly classified into 2 types.

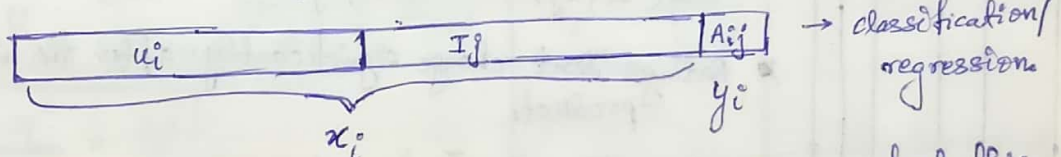
- ① Content based RS
- ② Collaborative filtering based RS.

## ① Content based RS

This requires feature vectors for user & items.



Now, we pose the training pt as



We won't discuss about this in this chapter.

### ③ collaborative filtering based RS

$$u_1 := M_{11}, M_{12}, M_{13}$$

$$u_2 := M_{21}, M_{23}, M_{24}$$

$$u_3 := M_{31}, M_{32}, M_{34}$$

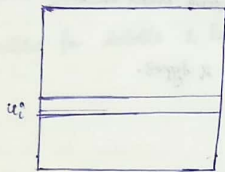
core-idea(assumption): users who agreed in the past tend to also agree in future.

We'll focus on this method in this chapter.

collaborative filtering based RS can further be divided into

- ① item-item based RS. (popularized by Amazon).
- ② user-user based RS.
- ③ Matrix factorization. (popular after Netflix 1 Million Challenge) (Matrix completion).

#### user-user based RS

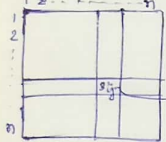


$u_i$  can be thought of as user-vector, very sparse like row.

So, we can find similarity b/w 2 users as

$$\text{sim}(u_i, u_j) = \cos(u_i, u_j) = \frac{u_i^T u_j}{\|u_i\| \|u_j\|}$$

$S_{ij}^u$  = user similarity matrix.



similarity( $u_i, u_j$ )

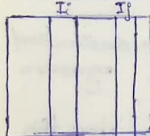
Tasks - recommend new items for  $u_{10}$ .

$u_{10} \rightarrow u_1, u_2, u_7$  are most similar to  $u_{10}$ .

\* Pick items liked by  $u_1, u_2, u_7$  that aren't yet watched by  $u_{10}$ .

Problem with user-user similarity  
user preference changes over time.

#### Item-Item based similarity



$$\text{sim}(I_i, I_j) = \cos(I_i, I_j)$$

\* Ratings don't change significantly after the initial product.

$u_{10} \rightarrow I_1, I_2, I_{10}$   
Find items similar to ones liked by user & recommends

#### Rule of Thumb $\rightarrow$

more users than items  $\rightarrow$  Amazon, Netflix, Youtube, etc.  
 $\downarrow$   
item ratings don't change much over time after initial period.  
 $\downarrow$   
prefer item-item over user-user.

#### Matrix factorization

$$6 = 2 \times 3 \text{ factors of } 6.$$

Similarly,  $A = BCD$  factors of  $A$ , and the process of breaking down  $A$  to  $B, C \times D$  is factorization/decomposition.

#### PCA as matrix factorization (Eigen Decomposition)

$X_{n \times d} \rightarrow$  data matrix (standardized).

$$S_{d \times d} = \frac{1}{n} (X^T X) = \text{co-variance matrix.}$$

{special type of fact<sup>n</sup> decomposition}

$$S_{d \times d} = W_{d \times d} \Lambda_{d \times d} W_{d \times d}^T \text{ (eigen decomposition of } S)$$

$$W = \begin{bmatrix} | & | & | & \dots & | \\ w_1 & w_2 & w_3 & \dots & w_d \\ | & | & | & \dots & | \end{bmatrix} \quad \Lambda = \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \lambda_3 & \\ & & & \ddots \end{bmatrix}$$

eigen values

So, from algebra math point of view,

PCA - eigen decomposition of co-variance matrix.

#### SVD (Singular Value Decomposition)

Another special type of decomposition

$$X_{n \times d} = U_{n \times n} \Sigma_{n \times d} V_{d \times d}^T$$

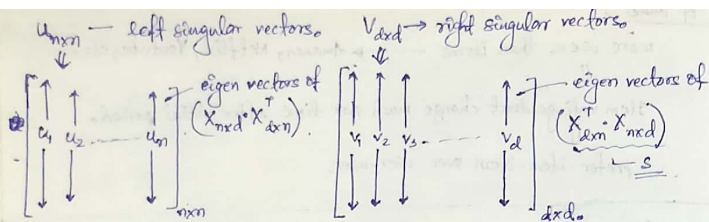
$$\Sigma_{n \times d} = \begin{bmatrix} s_1 & & & 0 \\ & s_2 & & \\ & & s_3 & \\ & & & \ddots \end{bmatrix}$$

diagonal matrix

singular values

$$\frac{s_i^2}{(n-1)} = \frac{\lambda_i^2}{(n-1)} \text{ eigen-values of } S_{d \times d}$$





### Non-negative Matrix factorization (NMF).

$$A_{n \times m} = B_{n \times d} \cdot C_{d \times m}$$

s.t.  $\begin{cases} B_{ij} \geq 0 \\ C_{ij} \geq 0 \end{cases}$  — constraint  
 We'll know later why it maybe useful.

### Matrix factorization for collaborative filtering

$A = \begin{bmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix}$   $A_{ij}$  — rating on  $I_j$  by  $u_i$   
 sparse matrix — many empty cells.

say, somehow we decompose  $A$  as:

$$A = B \cdot C^T \quad 0 < d \leq \min(m, n)$$

$B = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}$   $C = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}$   
 $B_i$   $C_j$

$A_{ij} = B_i \cdot C_j^T$  — assuming  $B_i$  is  $i$ th row of  $B$  & a column vector, &  $C_j$  is  $j$ th row of  $C$  & a column vector.

optimization problem —

solve using SGD  

$$\arg \min_{B, C} \sum_{i,j} (A_{ij} - B_i \cdot C_j^T)^2$$
 s.t.  $A_{ij}$  not empty.

solving this will give  $B$  &  $C$  as mentioned above  
 $\hat{A} = B \cdot C^T$  will give the predicted ratings

### Matrix factorization as feature engineering

Just in the above section,

$$A = B \cdot C^T$$

$B = n \times d$   
 $C = m \times d$  } these matrices can be thought of as featureizations of users & items.

similarly later we'll see word-vectors & eigen-faces — using MF.

### clustering as MF

In K-Means, the original optimization function was

$$\min_{c_i} \sum_{i=1}^k \sum_{x \in S_i} \|x - c_i\|^2 \quad x \in \mathbb{R}^d$$

We define  $x$  s.t.  $x_{ij} = \begin{cases} 1, & \text{if } x_j \in S_i \\ 0, & \text{otherwise.} \end{cases}$

$$Z = \begin{bmatrix} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{bmatrix}$$
  $Z_{ij}$

The above optimization function can be rewritten as

$$\min_{c_i, z_{ij}} \sum_{i=1}^k \sum_{j=1}^n z_{ij} \|x_j - c_i\|^2$$

if we have  $X = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}$  &  $C = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$   
 $B_i$   $C_j$

we can re-write the optimization function with vectorization as

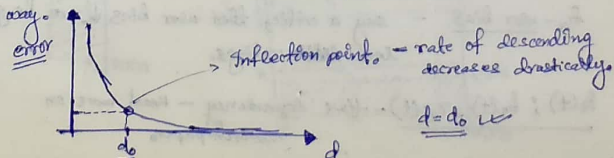
$$\min_{C, Z} \|X - CZ\|^2 \quad \text{s.t. } Z_{ij} = 0 \text{ or } 1 \quad \text{— NP-hard}$$

$\& \text{ column\_sum}(Z) = 1.$   
 So, clustering = MF + column-constraint + 0-1 constraint.

### Hyper-parameter tuning

① problem specific — depends on context & requirement.

② systematic way.



Netflix Prize (2009).

loss-metric - RMSE  $\sqrt{\frac{1}{n} \sum_{ij} (r_{ij} - \hat{r}_{ij})^2}$

Research paper → "datajobs.com/data-science-repo/Recommender-Systems-[Netflix].pdf".

$r_{ui}$  - rating by  $u$  to  $i$ .

$q_i$  - item vector.

$p_u$  - user vector.

Initial formulation

$$\min_{q_i, p_u} \sum_{ij} (r_{ij} - q_i^T p_u)^2 + \lambda (\|q_i\|^2 + \|p_u\|^2)$$

sq. loss.                      L2-regularization

Solves - SGD - find  $\frac{\partial L}{\partial p_u}$  &  $\frac{\partial L}{\partial q_i}$  & perform an iteration of GD.

It solves but takes long time.

② Alternate Least Squares (ALS).

fastest Algo. → ① fix  $p_u$ , gradient descent to fix  $q_i$ 's.  
→ ② fix  $q_i$ , gradient descent to fix  $p_u$ 's.

Modified formulation -

$$\min_{q_i, b_i, p_u, b_u} \sum_{ij} (r_{ij} - \mu - b_u - b_i - p_u^T q_i)^2 + \lambda \{ \|q_i\|^2 + \|p_u\|^2 + b_u^2 + b_i^2 \}$$

$\mu$  - avg rating across all users & items (say 3.0).  
It's better to separate the base avg fn interaction of user & items

$b_i$  - item bias - say a Nolan movie, so, it's supposed have a better rating. So, item bias ↑.

$b_u$  - user bias - say a critic, then user bias ↓, as he gives less ~~like~~ ratings.

$b_i(t)$ ;  $b_u(t)$ ;  $r_{ui}(t)$  - time dependency - Read more on research papers.

cold-start problem

↳ New user ( $u$ )  
or new product ( $i$ )

So, we have no rating data of them

New-user → Recommend top items based on meta-data.

\* geo-location (ip-address) - california } content based recommendations  
\* browser → safari  
\* device → ipad

new item → meta-data

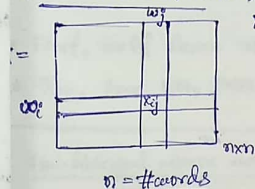
$i$  → category, price, descriptions. ....  
recommend as new arrival  $u_j$  →  $i_1, i_2, i_3$  ... similar items.

word-vectors using MF (truncated SVD)

word2vec → formed with help of neural networks.

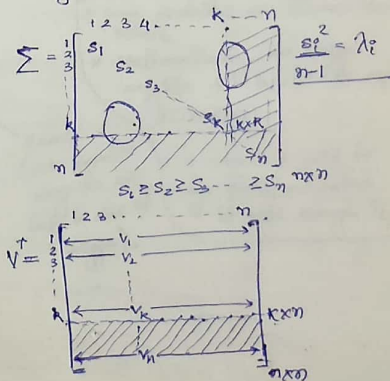
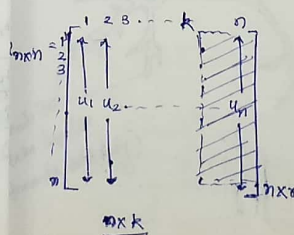
lots of other word vectors → LSA, PLSA, LDA, GLOVE (2014).  
related to word2vec making techniques.

① co-occurrence matrix ( $X$ )



$w_1, w_2, w_3, w_4, w_5, w_6$   
say two words are said to be in context of each other iff they occur within neighbourhood of distance 5.

$$X_{n \times n} = U_{n \times k} \Sigma_{k \times k} V_{k \times n}^T$$





PCA - top  $k$ -eigen-values & corresponding vectors. — analogous to what we did

$$X \approx \hat{X} = \hat{U} \sum_{k \times k} \hat{V}^T$$

$n \times n$       $n \times k$       $k \times k$       $k \times n$

$\hat{u}_i = i^{\text{th}}$  row of  $\hat{U}$  of  $k$ -dimensions can be thought of as the word vector of  $k$ -dimensions

$$u_i \in \mathbb{R}^k$$

problem with concurrence matrix ( $X$ )

↓

$n = \# \text{ words}$  very large.

sol<sup>n</sup>: use tf-idf to get top words, & build  $X$  from them.

eigen-faces is nothing but PCA on image data.

The eigen-vectors are basically called the eigen-faces.

