**Collaborative filtering:**

A hasn't watched All the movies in this case.

eg. $A = B \Rightarrow$ similar Ratings

$A + C : \Rightarrow$ const dissimilar cases

DATE

| Users | | HM1 | HM2 | HP3 | TW | SW | 2/3 |
|---|---|---|---|---|---|---|---|
| A | A | 4 | | | | 5 | 1 |
| | B | 5 | 5 | 4 | | | |
| | C | | | | 2 | 4 | 5 |
| | D | | 3 | | | | 3 |

movies →

**Jaccard Similarity.**

$$Sim(A,B) = \frac{2A \cap 2B}{2A \cup 2B}$$
$$= \frac{1}{5}$$

$$Sim(A,C) = \frac{2}{4}$$

$\therefore Sim(A,B) < Sim(A,C)$

just seen A, B have watched lesser movies than A, C have in common & not if they both liked it. problem

**Cosine Similarity**

→ Insert unknown values : 0.

$$Sim(A,B) = \frac{\sum A \cdot B}{\sqrt{A^2} \sqrt{B^2}} = \frac{4 * 5}{\sqrt{4^2 + 5^2 + 1^2} \cdot \sqrt{5^2 + 5^2 + 4^2}} = 0.38.$$

$Sim(A,C) = 0.32$ ← cosine ∠ b/w both.

$\therefore Sim(A,B) > Sim(A,C)$ but not by much.

problem: treats missing Ratings as least

**Centered Cosine** / **Pearson Correlation** / Rating Normaliz.

Normalization:

A : $(4 + 5 + 1) / 3 = 10/3$
B : $(5 + 5 + 4) / 3 = 14/3$
C : $(2 + 4 + 5) / 3 = 11/3$
D : $(3 + 3) / 2 = 6/2$

low Ratingbecomes -ve
High — +ve
then Subtract fromm Avg

Subtract this value form each cell.

ie. $4 - \frac{10}{3} = \frac{2}{3}$, $5 - \frac{10}{3} = $ ---

Can use Euclidean dist as well

| | HP1 | HP2 | HP3 | TW | SW1 | SW2 | SW3 | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 2/3 | | | 5/3 | -7/3 | | | | | | $\to \Sigma = 0$ | |
| B | 1/3 | 1/3 | -2/3 | | | | | | | | $\Sigma = 0$ | |
| C | | | | -5/3 | 1/3 | 4/3 | | | | | ✓ | |
| D | | 0 | | | | | 0 | | | | ✓ | |

this makes Avg. Rating ∀ users = 0 .

∴



**Before**
1) Avg = 2·5
2) Range = 0-5
3) 0 - 2·5 = Below Avg
4) Missing = 0

**After**
1) Avg = 0 .
2) -1 to +1
3) Below 0 : Below Avg.
4) Missing = Avg .



Now Compute Cosines.

$$Sim(A,B) = \frac{\frac{2}{3} * \frac{1}{3}}{\sqrt{\left(\frac{2}{3}\right)^2 + \left(\frac{5}{3}\right)^2 + \left(\frac{7}{3}\right)^2} \cdot \sqrt{\left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2 + \left(\frac{-2}{3}\right)^2}}$$

$$= \frac{2/9}{\sqrt{140/729}} \simeq 0.092$$

$Sim(A,C) = -0.56$ .

∴ $Sim(A,B) > Sim(A,C)$ .

<mark>ITEM ITEM</mark> Collaborative filtering
(More useful, outperforms user-user (CF)
∵ items are simpler than users
items : fixed set of genres .
↳ users : kinda like fuzzy logic Nature.

Because you watched Crime parrol, you might als
like sardhaan India .

precision @ k , Recall @ k
↓
#K Recommn".
#Actual Relevant

# estimate Rating done by user5 for movie1

Left margin:
=3.6
3·16
?
·4
3
·6

**Ratings table** (movies × users)

| movies \ users | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | | 3 | 4 | ? | 5 | | | 5 | | 4 | |
| 2 | | | | | | | | | | 2 | 1 | 3 |
| 3 | 2 | 4 | 5 | 1 | 2 | 5 | 3 | 4 | 4 | 3 | 5 | |
| 4 | | 2 | 4 | | 5 | 4 | | 4 | | | 2 | |
| 5 | | | 4 | 3 | 4 | 2 | | | | | 2 | 5 |
| 6 | 1 | | 3 | | 3 | | | 2 | | | 4 | |

**Mean-centered table**

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | -2.6 | 0 | -0.6 | 0 | ? | 1.4 | 0 | 0 | 1.4 | 0 | 0.4 | 0 |
| 2 | 0 | 0 | 1.84 | 0.84 | 0 | 0.84 | 0 | 0 | 0 | -1.16 | -2.16 | -0.16 |
| 3 | -1 | 1 | 0. | -2 | -1 | 0 | 0 | 0 | 1 | 0 | 2 | 0 |
| 4 | 0 | -1.4 | 0.6 | 0 | 1.6 | 0 | 0 | 0.6 | 0 | 0 | -1.4 | 0 |
| 5 | 0 | 0 | 0.7 | -0.3 | 0.7 | -1.3 | 0 | 0 | 0 | 0 | -1.3 | 1.7 |
| 6 | -1.6 | 0 | 0.4 | 0 | 0.4 | 0 | 0 | -0.6 | 0 | 0 | 1.4 | 0 |

$\text{sim}(M1, \text{other Movies Rated by user 5}) = ?$   2 movies highest similarity

$\text{sim}(m1, M3)$   $(1,5) = ?$

$$(1,3) = \frac{2.6*1 + 1.4*1 + 0.4*3}{\sqrt{-u}\;\;\sqrt{-u}} = 0.41 \;\checkmark$$

↑ take weighted Avg using these

$\text{sim}(m1, m4) = -0.10$

$\text{sim}(m1, m5) = -0.31$

$\text{sim}(m1, m6) = (1,6) = \underline{\qquad} = 0.59$

$\text{sim}(m1, m1) = 1$

$\text{sim}(m1, m2) = -0.18$

$(3,5)\qquad (6,5)$

$$\therefore \text{Ratings for } \boxed{?} \Rightarrow \frac{0.41*2 + 0.59*3}{0.41+0.59} = 2.6$$

$$\therefore \boxed{?} = 2.6$$

④

User preferences $\longrightarrow$ Recommender System $\longrightarrow$ Recommendations (predicting future behaviour).

Explicit feedback
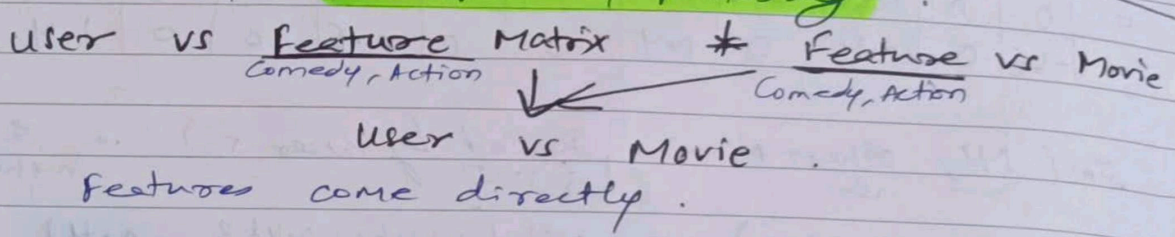
implicit feedback

**Collaborative:**

Similar users like Similar things.
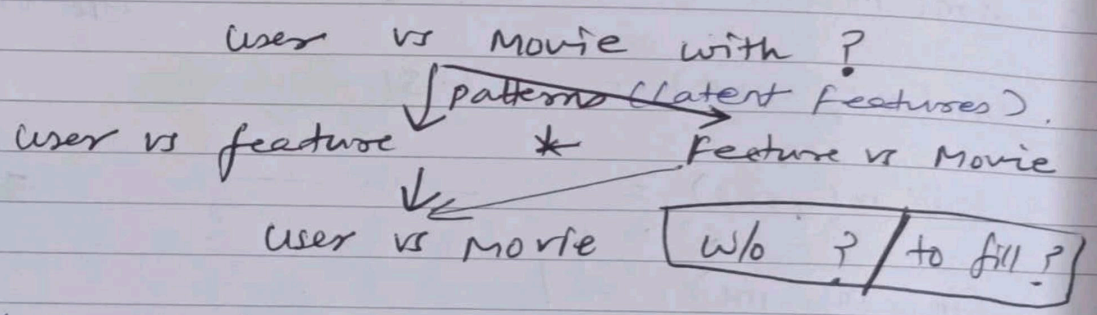user × item Matrix

**Content Based:**

Considers item/user feature

item × item
movie genre,
Years of Release,
Cast,
director,
prod: house.

User × User
age, gender,
Spoken language

---

1. ==Content Based Filtering==.

User vs **Feature** Matrix   *   **Feature** vs Movie
   Comedy, Action               Comedy, Action

       ↓

User vs Movie.
Features come directly.

2. ==Collaborative Based Filtering==.

User vs Movie with ?
      patterns (latent features).

User vs feature   *   Feature vs Movie

      ↓

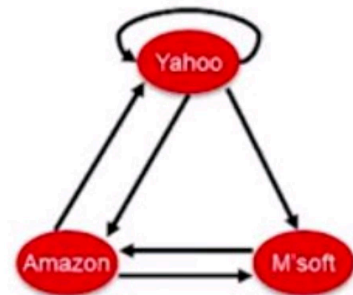User vs movie   | w/o ? / to fill ? |

you keep guessing values for both Matrixs until you
get closest Matrix Multipl^n product same as from where u
started.

Applications of Recommend? system : News/Songs/---
eg. Synthetic Control. what'll be effect of "Gun Control"
policies if Implemented? you check for countries → to
you that already Implemented em

# 8. Hyperlink Induced Topic Search (HITS) Algorithm

## Example

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \qquad A^T = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| h(yahoo) | = | .58 | .80 | .80 | .79 | $\cdots$ | .788 |
| h(amazon) | = | .58 | .53 | .53 | .57 | $\cdots$ | .577 |
| h(m'soft) | = | .58 | .27 | .27 | .23 | $\cdots$ | .211 |
| a(yahoo) | = | .58 | .58 | .62 | .62 | $\cdots$ | .628 |
| a(amazon) | = | .58 | .58 | .49 | .49 | $\cdots$ | .459 |
| a(m'soft) | = | .58 | .58 | .62 | .62 | $\cdots$ | .628 |

## 1. Content Based

```
        user profile              item profile
            ↓                          ↓
```

what users likes
└ user u likes MI
∴ user features like
  age, gender, job mattered ...
· User x feature Matrix

Features of items : #
genre, director actors ...
       feature * movie matrix
       ┌ like, dont like
       └ 0, 1, 2, 3, 5 scale of liking

utility matrix : User x movie

prefers

user ___ a good camera more
than the phones RAM

user x features x weight

| user | likes_features | w weight of |
|------|----------------|-------------|
| 1    | 1              | 2           |
|      | 3              | 1           |
|      | 4              | 1           |
| 2    |                |             |
| :    |                |             |
| :    |                |             |

Oppo phone has a
great camera.

product x features
+ weight

| prod1 | has_features | order of quality |
|-------|--------------|------------------|
| 1     | 1            | 1                |
|       | 3            | 1                |
|       | 4            | 1                |
| 2     | 2            | 1                |
|       | 3            | 4                |
| 3     | 1            | 3                |
|       | 4            | 1                |

| | feature1 | feature2 | feature3 | F4 |
|---|---|---|---|---|
| product1 | $\boxed{1}$ | | 1 | 1 |
| product2 | | 1 | 4 | |
| product3 | 3 | | | |
| | | | | |
| user1 | $\boxed{2}$ | | 1 | 1 |

∴ user Interest in product1 = $2*1 + 1*1 + 1* 1 = 5$

$2 = 1* 4 = 4$

$3 = 2\times3 + (1) = 7$

| user | prod | weight |
|------|------|--------|
| 1    | 1    | 5      |
|      | 2    | 4      |
|      | 3    | 7      |
| 2    | :    | :      |

$$\therefore (1,3) = \frac{(1,2) * (2,3) + (1,4) * (4,3)}{(1,2) + (1,4)}$$
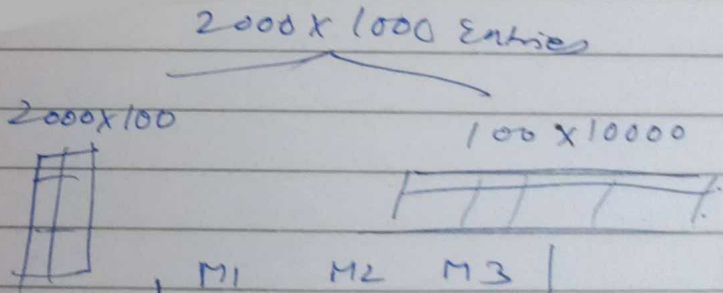
if sparsity too high, u cant figure out the relations

Movie.

user

|   | 1 | 2 | 3 |
|---|---|---|---|
| 1 | ④ | 3 | 2 |
| 2 | 1 | 3 | 3 |
| 3 | 5 | ④ 5 | |

Ratings $\Rightarrow$ Dependent on Rows & Cols

1) Pref of persons 3 = $\dfrac{\text{pref of persons } 1 + 2}{2}$

Comedy   Action

or its possible that

$\dfrac{M1 + M3}{2} = M2$

u try to approximate the original uxi matrix
from uxf and mxf matrix
that were randomly created to fill the ?
from the original uxi matrix

random uxf matrix   random fxm matrix

The Features are not Comedy/Action explicit.

|    | f1  | f2  |
|----|-----|-----|
| u1 | 0.2 | 0.5 |
| u2 | 0.3 | 0.4 |
| u3 | 0.7 | 0.8 |

3x2

2000x100

$2000 \times 1000$ Entries

$100 \times 10000$

|    | M1  | M2  | M3  |
|----|-----|-----|-----|
| f1 | 1.2 | 3.1 | 0.3 |
| f2 | 2.4 | 1.5 | 4.4 |

2x3

M1 M2 M3

| u1 |
| u2 |
| u3 |

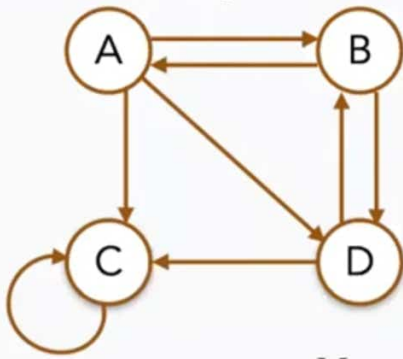$(u_1 f_1), (f_1 m_1)$
to $f(u_1 m_1)$

# Dead ends

## A tiny web



- ❖ Let's make C a dead end
- ❖ *M* is not stochastic anymore, rather *substochastic*
  - ▸ The 3rd column sum = 0 (not 1)
- ❖ Now the iteration $v := Mv$ takes all probabilities to zero

| $M$ | | | |
|---|---|---|---|
| 0 | 1/2 | **0** | 0 |
| 1/3 | 0 | **0** | 1/2 |
| 1/3 | 0 | **0** | 1/2 |
| 1/3 | 1/2 | **0** | 0 |

| $v$ |
|---|
| 1/4 |
| 1/4 |
| 1/4 |
| 1/4 |

| $Mv$ |
|---|
| 3/24 |
| 5/24 |
| 5/24 |
| 5/24 |

| $M^2v$ |
|---|
| 5/48 |
| 7/48 |
| 7/48 |
| 7/48 |

$\cdots \longrightarrow$

| |
|---|
| 0 |
| 0 |
| 0 |
| 0 |

# Spider traps

## A tiny web



- Let C be a one node spider trap
- Now the iteration $v := Mv$ takes all probabilities to zero except the spider trap
- The spider trap gets all the PageRank

| | $M$ | | |
|---|---|---|---|
| 0 | 1/2 | 0 | 0 |
| 1/3 | 0 | 0 | 1/2 |
| 1/3 | 0 | 1 | 1/2 |
| 1/3 | 1/2 | 0 | 0 |

| $v$ |
|---|
| 1/4 |
| 1/4 |
| 1/4 |
| 1/4 |

| $Mv$ |
|---|
| 3/24 |
| 5/24 |
| 11/24 |
| 5/24 |

| $M^2v$ |
|---|
| 5/48 |
| 7/48 |
| 29/48 |
| 7/48 |

$\cdots \longrightarrow$

| |
|---|
| 0 |
| 0 |
| 1 |
| 0 |

# Some Problems with Page Rank

- **Measures generic popularity of a page**
  - Biased against topic-specific authorities
  - **Solution:** Topic-Specific PageRank (**next**)
- **Uses a single measure of importance**
  - Other models of importance
  - **Solution:** Hubs-and-Authorities
- **Susceptible to Link spam**
  - Artificial link topographies created in order to boost page rank
  - **Solution:** TrustRank

In this section we shall discuss about the following distance measures in details:

> (1) Euclidean Distance
>
> (2) Jaccard Distance
>
> (3) Cosine Distance
>
> (4) Edit Distance
>
> (5) Hamming Distance

## 6.1.1 Euclidean Distances

**Q. What do you mean by Euclidean distance ? Explain with example.**

- The Euclidean distance is the most popular out of all the different distance measures

- The Euclidean distance is measured on the Euclidean space. If we consider an n-dimensional then each point in that space is a vector of n real numbers. For example, if we consider Euclidean space then each point in the space is represented by $(x_1, x_2)$ where $x_1$ and $x_2$ are

- The most familiar Euclidean distance measure is known as the $L_2$– norm which in the defined as :

$$d([x_1, x_2,...x_n], [y_1, y_2,...y_n]) = \sqrt{\sum_{i=1}^{n} (x_i - y_i)^2}$$

or the two-dimensional space the $L_2$– norm will be :

$$d([x_1, x_2], [y_1, y_2]) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

- We can easily verify all the distance axioms on the Euclidean distance :

## 6 | Real Time Big Data Model

### Iteration II:

$$A^2 \cdot V = A \cdot (A^1 \cdot V)$$

$$\begin{bmatrix} 0.4 \\ 0.2 \\ 0.2 \\ 0.05 \\ 0.05 \\ 0.1 \end{bmatrix}$$