

→ The heart of collaborative filtering is the ability to find people similar to you.

Ways to measure Similarity:

1. Cosine Similarity \Rightarrow used to measure how similar two items are

Items are user behaviour data.

$$\text{Cossim}(x, y) = \frac{\sum_i x_i y_i}{\sqrt{\sum_i x_i^2} \sqrt{\sum_i y_i^2}}$$

In matrix form each attribute is a dimension. So if have n attributes it forms n -dimensional space.

↳ angle b/w two attributes is the similarity cos of those two attributes.

(2) Sparsity:
These attributes are for user behaviour only not for content.

Data is sparse (small amount of data is spreaded over a large area)

↳ Because, there are many movies in the world and no one sees all the movies.

↳ ~~So some~~

↳ So, one person who has seen a particular movie; ~~then~~ and that particular movie is may/may not be seen by other persons.

↳ To find similar person who watched the movie is difficult. So, it's tough for Collaborative Filtering to find similar person.

↳ This collaborative filter is working well for Amazon, Netflix because they have more no. of users; so they will generate meaningful relations.

Similarity Metrics:-

① Adjusted Cosine:-

$$\text{CosSim}(x, y) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2} \sqrt{\sum_i (y_i - \bar{y})^2}}$$

ratings are different for opposition to ~~oppn~~

↳ Some countries are more brutal with their ratings than others.

↳ so, adjusted cosines attempts to normalize these differences.

↳ If we have less data (people who rated the movie are very less) then we don't get a meaningful average.

↳ so, the difference b/w user & users mean is "0"

↳ If we have large data & many people have rated means then Adjusted cosine is useful.

(2) (Item-based) Pearson Similarity:

$$\text{CosSim}(x, y) = \frac{\sum_i (x_i - \bar{i})(y_i - \bar{i})}{\sqrt{\sum_i (x_i - \bar{i})^2} \sqrt{\sum_i (y_i - \bar{i})^2}}$$

Here we look at difference b/w ratings & the average from all users for that given item.

\bar{i} → average ~~is~~ rating of the item.

↳ If we have sparse data, it is a good metric.

③ Spearman rank Correlation:-

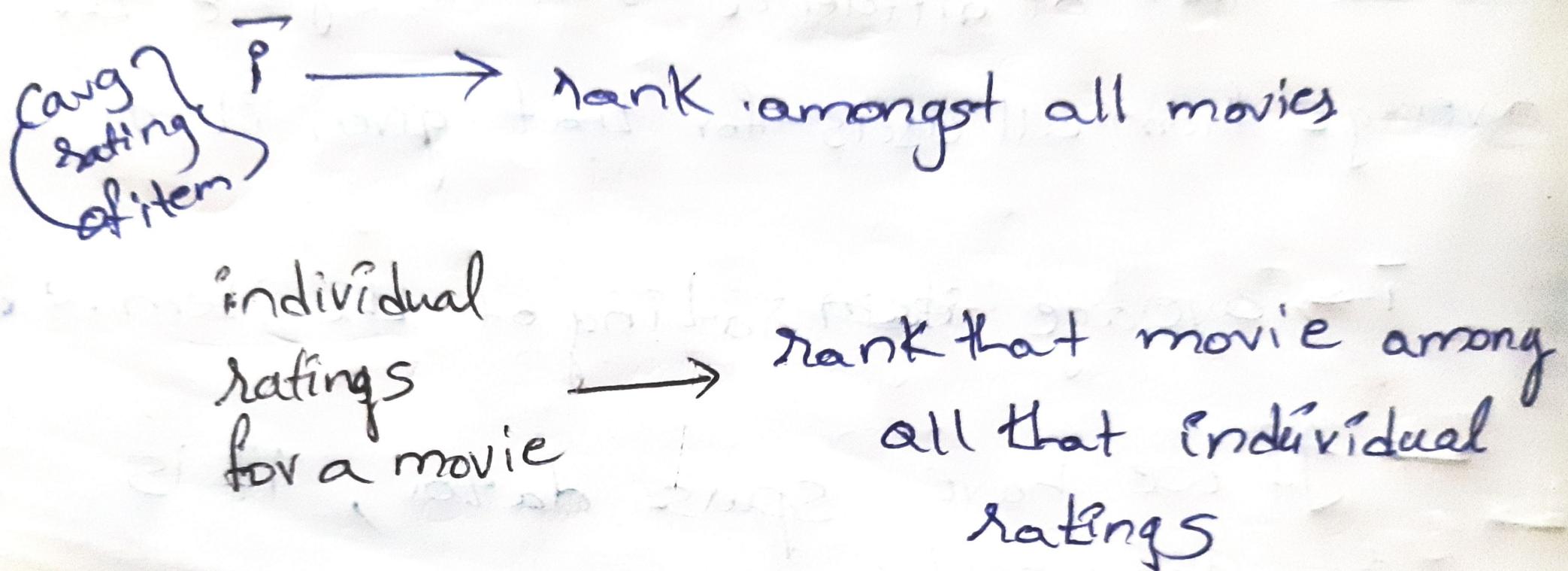
Pearson Similarity ~~rankings~~
based on ranks

not ratings

↳ same ideas as Pearson's Similarity

but instead of ratings scores directly

we use ranks instead.



(ii) Mean squared Difference

$$MSD(x, y) = \frac{\sum_{i \in I_{xy}} (x_i - y_i)^2}{|I_{xy}|}$$

$$MSD_{Sim}(x, y) = \frac{1}{MSD(x, y) + 1}$$

→ Here you just directly comparing how two people rated the same set of things.

→ ~~Its very much same idea~~

ITS Same idea as mean absolute error

$|I_{xy}| \rightarrow$ no. of items each user had

in common that we summed across to get the average/mean.

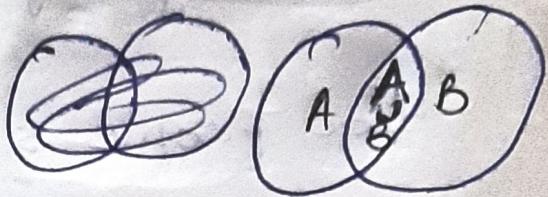
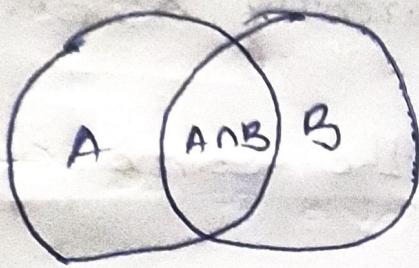
→ $MSD(x, y)$
The problem is, we computed how different users $x \& y$ are

→ To get similarity we will take inverse

i.e So, $MSD_{Sim}(x, y) + 1$

$X \& Y \rightarrow$ two diff things; not diff people in order to avoid 0

⑤ Jaccard Similarity:



$A \& B \rightarrow$ two users.

$$= \frac{A \cap B \text{ ratings}}{A \cup B \text{ ratings}}$$

⇒ no. of user A ratings
intersection with
no. of user B ratings

as we are counting things & we
are not looking the actual ratings

values. So, we are throwing some
information

→ It is used for implicit ratings

because we know their interests but
by taking their ratings (so, ratings won't come)