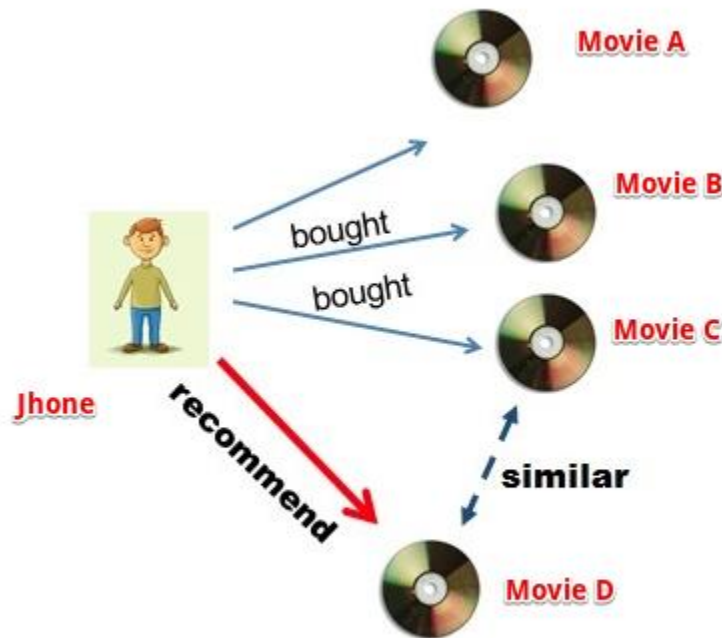


Item Based Collaborative Filtering

Item-item collaborative filtering is one kind of recommendation method which looks for similar items based on the items users have already liked or positively interacted with. It was developed by Amazon in 1998 and plays a great role in Amazon's success.



How IBCF works are that it suggests an item based on items the user has previously consumed. It looks for the items the user has consumed then it finds other items similar to consumed items and recommends accordingly.

Let's understand this with an example. Suppose our user Jhone wants to purchase a movie DVD. Our job is to recommend him a movie based on his past preferences. We will first search for movies that Jhone has watched or liked, let's call those movies 'A', 'B' and 'C'. Next, we will search for other movies similar to three movies. Suppose we found out that movie 'D' is highly similar to 'C', therefore, there is a highly likely chance that Jhone will also like movie 'D' because it is similar to one Jhone has already liked. Hence, we will suggest the movie 'D' to Jhone.

So at its core IBCF is all about finding items similar to the ones user has already liked. But how to find similar items? and what if there are multiple similar items in that case which item to suggest first? To understand this lets' first understand the intuition behind the process, this will assist us to comprehend the mathematics behind the IBCF recommendation filtering.

Finding Similarity Between Items

Suppose we have a table of users and their ratings for movies

User	1: Toy Story	2: Star Wars: Epi	356: Forrest Gump	318: Shawshank Redemption	593: Silence of the Lambs	3578: Gladiator
755	2	5	2		4	4
5277	1			2	4	2
1577				5	2	
4388	2	3				1
1202		3	4	1	4	1
3823	3	4	4	4		
5448			3	1	1	4
5347	2				3	2
4117	4	1		4	2	4
2765	4	2		5	3	
5450	5	1	5			5
139	2	5	2			
1940	4			5	4	
3118	3		3		2	
4656	2	4			5	5

Let's pick two movies (Id 1: Toy Story and Id 2: Star Wars) for which we have to calculate the similarity i.e. how much these two movies are comparable to one another in terms of their likeness by users. To compute this we will:

First multiple ratings of both movies with each other and the sum of the result.
Let's call this value 'A'.

User	1: Toy Story	2: Star Wars	Product of Movie Ratings	
755	2	5	10	
5277	1		0	
1577			0	
4388	2	3	6	
1202		3	0	
3823	3	4	12	
5448			0	
5347	2		0	
4117	4	1	4	← Product of ratings
2765	4	2	8	
5450	5	1	5	
139	2	5	10	
1940	4		0	
3118	3		0	
4656	2	4	8	
4796			0	
6037	2		0	
3048	4	5	20	
4790	2	1	2	
4489	2	2	4	
		A = 89		← Sum of product of ratings

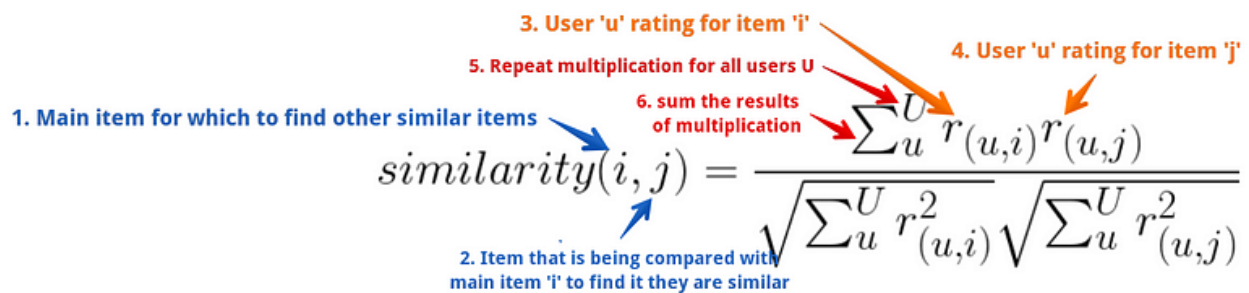
similarity between movie 1 and 2

Repeating the above process for all the movies will result in a table with similarities between each movie(in general we call them items).

Here is how the above process is depicted in mathematical form.

$$similarity(i, j) = \frac{\sum_u^U r(u, i) r(u, j)}{\sqrt{\sum_u^U r(u, i)^2} \sqrt{\sum_u^U r(u, j)^2}}$$

Adding some labels to the letters will ease in understanding of each part of the equation.



The diagram shows the similarity equation with six numbered annotations and arrows pointing to specific parts of the formula:

- 1. Main item for which to find other similar items**: Points to the variable i in the numerator.
- 2. Item that is being compared with main item 'i' to find if they are similar**: Points to the variable j in the numerator.
- 3. User 'u' rating for item 'i'**: Points to the term $r(u, i)$ in the numerator.
- 4. User 'u' rating for item 'j'**: Points to the term $r(u, j)$ in the numerator.
- 5. Repeat multiplication for all users U**: Points to the summation symbol \sum_u^U in the numerator.
- 6. sum the results of multiplication**: Points to the summation symbol \sum_u^U in the denominator.

$$similarity(i, j) = \frac{\sum_u^U r(u, i) r(u, j)}{\sqrt{\sum_u^U r(u, i)^2} \sqrt{\sum_u^U r(u, j)^2}}$$

