

## ★ Cache() & Persist() in spark →

### ~~Cache()~~ →

- Stores in-memory only.
- If dataset is reused cache it
- in Collaborative filtering, similarity matrix is queried many times  
→ caching avoid recomputation

### Cache(i)

- Stores in memory only (RAM)
- good for small dataset which fit in RAM.
- shortcut of .persist(StorageLevel.MEMORY-ONLY)

### persist()

- More flexible : Can store to disk, memory & both.
- useful for data too big for RAM.
- .persist(StorageLevel.MEMORY-AND-DISK)

- ★ Item Based Collaborative filtering.  
Recommendation technique that predict a user's preferences for an item based on their rating of similar items

Analogy → Sia Rated Fight Club & Forrest Gump similarly  
So, do Ben,  
Ted only watched Forrest Gump, so we can recommend Fight Club to him.  
very old school technique.

Steps →

1. Self-Join on userid → get (movie1, movie2, rating1, rating2)
2. Filter out duplicates (movie1, movie2) ≠ (movie2, movie1)
3. Compute cosine similarity

$$\cos(A, B) = \frac{\sum (r_A \cdot r_B)}{\sqrt{\sum (r_A^2)} \cdot \sqrt{\sum (r_B^2)}}$$

4. Rank & Pick Top similar movies.

why? →  
because Captures pattern of similarity even if user's rate on different scales. (Normalizes Ratings)

- ★ System module in python →

```
import sys → captures cli arguments  
if (len(sys.argv) > 1) → check if movie id is passed or not  
we can → spark-submit file.py 50  
sys.exit(1) & sys.exit(0)
```

★ Code for movie similarity  $\rightarrow$

1. Core logic  $\rightarrow$  Cosine similarity (def Cosine Compute Similarity(date):)

$$xx = \text{rating}_1^2, \quad yy = \text{rating}_2^2, \quad xy = \text{rating}_1 * \text{rating}_2$$

group by (movie1, movie2)

$$e) \frac{\text{numerator}}{\text{denominator}} = \frac{\sum(xy)}{\sqrt{\sum(xx)} * \sqrt{\sum(yy)}} = \text{Score}$$

if denominator  $\neq 0 \rightarrow \text{score} = 0$

output  $\rightarrow$  (movie 1, movie 2, score, numpairs)  $\rightarrow$  count(xy)

2. after return of dataframe lacks it.

3. Threshold for good similarity  $\rightarrow$

- Score Threshold = 0.97  $\rightarrow$  only recommend if similarity is High
- coOccurrence Threshold = 50, if number of pairs (enough users) rated both.

4. Broadcast lookup  $\rightarrow$   $(\text{def } \text{getMovieName} (\text{movies name}, \text{movie ID}))$   
uses movieName dataset to Map IDs  $\rightarrow$  titles.

Short flow  $\rightarrow$

load  $\rightarrow$  Pair  $\rightarrow$  Cosine  $\rightarrow$  Cache  $\rightarrow$  Filter  $\rightarrow$  Sort  $\rightarrow$  Print.