

TRƯỜNG ĐẠI HỌC CÔNG NGHỆ - ĐẠI HỌC QUỐC GIA HÀ NỘI  
VIỆN TRÍ TUỆ NHÂN TẠO

---

BÁO CÁO MÔN HỌC KỸ THUẬT VÀ CÔNG NGHỆ DỮ LIỆU LỚN

ĐỀ TÀI

Phân tích dữ liệu lớn  
đề xuất sản phẩm cho người dùng

Nhóm sinh viên thực hiện: Đỗ Minh Nhật - Trần Nam Anh - Hà Như Ý

Giảng viên hướng dẫn: TS. Trần Hồng Việt

Ths. Ngô Minh Hương

Hà Nội, ngày 27 tháng 05 năm 2024



## Chương 2: Thuật toán gợi ý sản phẩm

### 2.1. Thuật toán content-based filtering

#### 2.1.1. Giới thiệu

Hệ thống này đề xuất nội dung cho người dùng dựa trên nội dung mà người dùng đó thích. Ví dụ, người dùng A thích xem bộ phim có thể loại hài hước, vui vẻ, trong tương lai hệ thống sẽ đề xuất các nội dung liên quan đến bộ phim có thể loại hài hước, vui vẻ. Cách tiếp cận này yêu cầu phải có vector đặc trưng của mỗi item, không chỉ vậy đánh giá của mỗi người dùng và item cũng rất quan trọng vì nếu không có user rating của user đối với item ta không thể nào biết được user đó có sở thích là gì để có thể đề xuất nội dung tương tự. Tuy nhiên, đôi lúc trong thực tế việc xác định feature vector của item có thể gặp khó khăn do khó xác định được nhóm cụ thể các item.

#### 2.1.2. Method.

##### 2.1.2.1 Xây dựng Utility Matrix.

Việc xây dựng Utility Matrix thường rất quan trọng, nếu không có Utility Matrix dường như ta không thể gợi ý được sản phẩm tới người dùng, ngoại trừ cách gợi ý các sản phẩm phổ biến nhất. Do đó, việc xây dựng Utility Matrix trong bài toán recommendation là cực kỳ quan trọng. Có 2 thực thể chính trong Recommendation System đó là users và items. Trong đó, mỗi user sẽ có mức độ quan tâm tới từng item. Trong thực tế, mức độ quan tâm đó thường được phân tích dựa trên rating của user đối với item, hoặc like/dislike của user đối với item đó.

	A	B	C	D	E	F
Mưa nửa đêm	5	5	0	0	1	?
Cỏ úa	5	?	?	0	?	?
Vùng lá me bay	?	4	1	?	?	1
Con cò bé bé	1	1	4	4	4	?
Em yêu trường em	1	0	5	?	?	?

Hình 1: Ví dụ về utility matrix với hệ thống Gợi ý bài hát. Các bài hát được người dùng đánh giá theo mức độ từ 0 đến 5 sao. Các dấu '?' nền màu xám ứng với việc dữ liệu chưa tồn tại trong cơ sở dữ liệu. Recommendation Systems cần phải *tự điền* các giá trị này.

**Hình 1:** Utility Matrix mô tả mức độ yêu thích của user với item.

**Nguồn:** <https://machinelearningcoban.com/2017/05/17/contentbasedrecommendersys/>

##### 2.1.2.2 Xây dựng Item Profiles.

Trong Content-based Recommendation System, hệ thống đề xuất dựa trên nội dung. Chúng ta cần xây dựng một bộ hồ sơ cho mỗi item còn được gọi là profile item. Profile item này thường được biểu diễn dưới dạng vector thường được gọi là feature vector. Trong trường hợp đơn giản, item profile được trích xuất từ item. Ví dụ: một bộ phim khi sản xuất sẽ thường bao gồm các nội dung như: Đạo diễn, Năm phát hành, Thể loại,... từ thông tin đó ta có thể xây dựng feature vector cho từng bộ phim.

	A	B	C	D	E	F	item's feature vectors
Mưa nửa đêm	5	5	0	0	1	?	$\mathbf{x}_1 = [0.99, 0.02]$
Cỏ úa	5	?	?	0	?	?	$\mathbf{x}_2 = [0.91, 0.11]$
Vùng lá me bay	?	4	1	?	?	1	$\mathbf{x}_3 = [0.95, 0.05]$
Con cò bé bé	1	1	4	4	4	?	$\mathbf{x}_4 = [0.01, 0.99]$
Em yêu trường em	1	0	5	?	?	?	$\mathbf{x}_5 = [0.03, 0.98]$
User's models	$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$\theta_5$	$\theta_6$	$\leftarrow$ need to optimize

**Hình 2:** Feature vector của bài hát được đặt ở cuối cùng. Các entry trong utility matrix  $(i, j)$  đại diện cho rating của item thứ  $i$  được user thứ  $j$  rated, trong đó  $(i, j)$  là vị trí hàng thứ  $i$  cột thứ  $j$ . Nguồn:

<https://machinelearningcoban.com/2017/05/17/contentbasedrecommendersys/>

Ví dụ ở hình 2 phía trên, ta xây dựng feature vector  $x_i$  với mỗi bài hát. Trong đó, chiều thứ nhất của vector  $x_i$  là số thực đo mức độ bài hát liên quan đến Bolero và chiều thứ 2 của vector  $x_i$  là số thực đo mức độ bài hát liên quan đến Thiếu nhi. Sau khi xây dựng trong feature vector, ta có thể chỉ ra rằng trong hình 1 user A thích nghe nhạc Bolero và không thích nghe nhạc thiếu nhi. User C thích nghe nhạc thiếu nhi và không thích nghe nhạc Bolero. Khi đó hệ thống sẽ ưu tiên đề xuất cho user A nhạc Bolero, còn user B sẽ được hệ thống đề xuất nhạc thiếu nhi.

### 2.1.2.3 Xây dựng Objective Function.

Với mỗi user trong toàn bộ user ta cần tìm một bộ tham số  $\theta_i$  bao gồm  $(w_i, b_i)$  sao cho ta có thể xấp xỉ hàm mục tiêu  $y$ . Gọi hàm xấp xỉ mục tiêu là  $\hat{y}_{mn}$ , ta có:

$$\hat{y}_{mn} = x_m w_n + b_n$$

Với  $x_i$  là vector hàng,  $w_j$  là vector cột.

Ở đây ta có thể xây dựng loss function tùy ý. Một trong các cách đơn giản thường được sử dụng đó là Square Error.

$$L_n = \frac{1}{2} \sum_{m: r_{mn}=1} (x_m w_n + b_n - y_{mn})^2$$

Trong đó  $x_m$  là feature vector của item,  $w_n$  là weight matrix đại diện cho sở thích của user với từng từng feature trong feature vector,  $b_n$  là bias vector,  $y_{mn}$  là target value ( đại diện cho rating của item thứ m được user thứ n rated ).  $\{m : r_{mn} = 1\}$  là tập các vị trí (m, n) trong matrix mà user thứ n đã rating cho item thứ m.

#### 2.1.2.4 Phương pháp tiếp cận sử dụng MapReduce

**Bước 1: Tính toán user profile** ( sở thích của user, ví dụ: xem phim hài hước, vui nhộn)

- *Mapper*: ( key, value ) => (user\_id, features\_of\_item\_id\_have\_interacted) (ở đây đối với item mà người dùng thích feature vector là chính nó, đối với item mà người dùng không thích feature vector sẽ bằng với giá trị âm của feature vector đó feature vector = -feature\_vector ).
- *Reducer* : (key, value ) => (user\_id, sum( features\_of\_item\_id\_have\_interacted )).

**Bước 2: Tính cosine similarity giữa user và movie:**

- *Mapper*: (key, value) => (user\_id, (movie\_id, score)).
- *Reducer*: Chọn một thresh hold có giá trị bằng A, nếu giá trị của score >= A giữ lại các cặp ( user\_id, (movie\_id, score)).

**Bước 3: Lọc các cặp có giá trị score nhỏ hơn threshold, giữ lại các cặp lớn hơn threshold tương ứng với các cặp có thể đề xuất kết bạn**

- *Mapper*: Truyền output của Reducer trong *Bước 2* cho Reducer.
- *Reducer*: Reduce về dạng (key,value) = (user\_id, movie\_id). Khi đó user\_id sẽ được đề xuất xem các movie\_id.

## 2.2 Collaborative Filtering

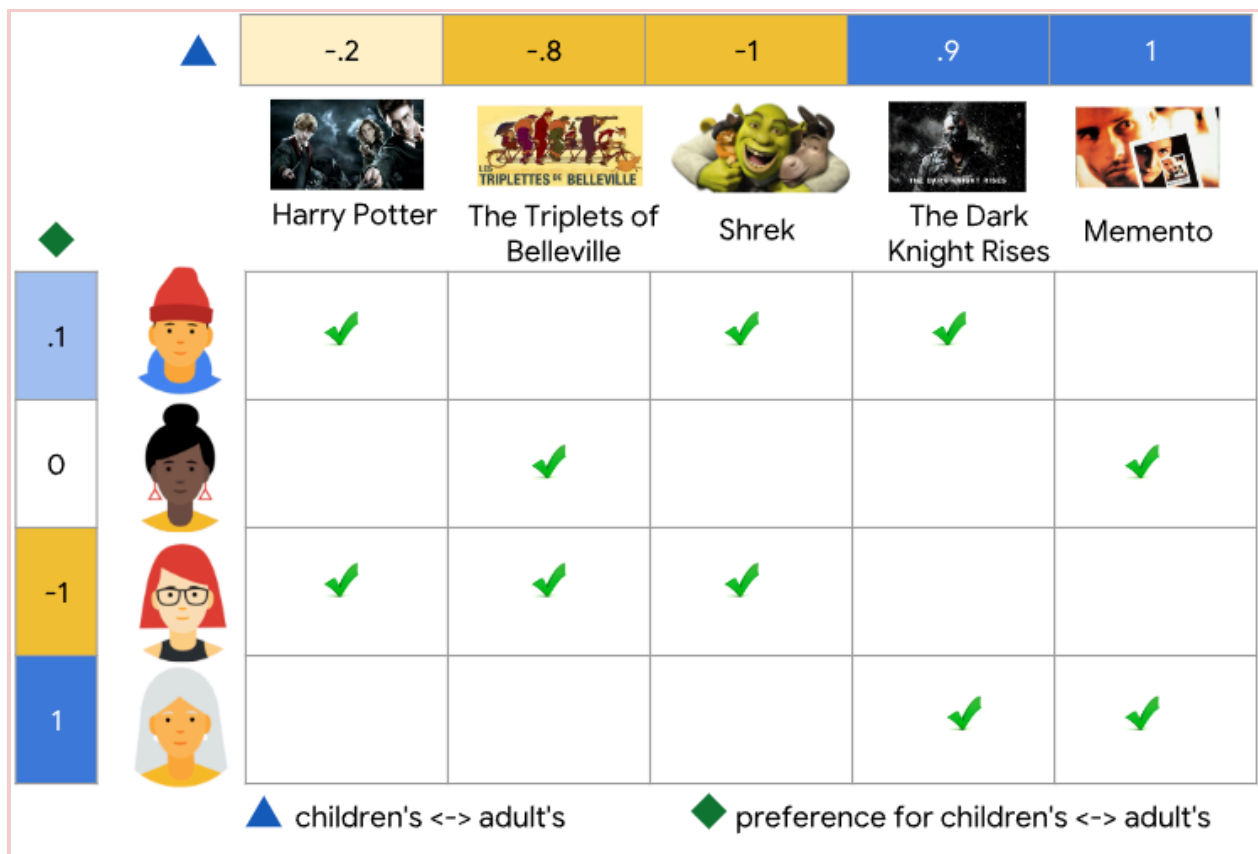
### 2.2.1 Introduction

Trong Content-based Recommendation System, chúng ta đã biết cách tạo ra một hệ thống đề xuất dựa trên các đặc trưng của mỗi item. Đặc trưng của Content-based Recommendation System là đề xuất dựa trên đặc trưng của mỗi item, nhưng câu hỏi đặt ra là nếu giả sử user không rating bất cứ item nào hoặc rất ít làm sao ta có thể đề xuất cho user tới item phù hợp. Trong thực tế, việc user rating item là rất ít, không phải người nào xem youtube cũng bấm like/dislike, không phải người nào sau khi mua hàng cũng đánh giá sản phẩm. Không chỉ vậy, trong thực tế có rất nhiều các user có thể có các sở thích giống nhau ví dụ như user A thích item X, item Y, item Z, user B thích item X, item Y, vậy có khả năng rất cao là user B cũng thích item Z. Câu hỏi đặt ra là:

1. Làm thế nào ta có thể xác định được sự giống nhau giữa các cặp users.
2. Khi đã biết các cặp users gần giống nhau làm thế nào để dự đoán mức độ quan tâm của user này dựa trên user khác mà gần giống với user này

### 2.2.2 Method

**2.2.2.1. 1D Embedding:** Giả sử chúng tôi gán cho mỗi bộ phim một số thực trong khoảng  $[-1, 1]$  mô tả bộ phim đó cho trẻ em hay cho người lớn với giá trị càng gần với -1 đại diện cho bộ phim dành cho trẻ em và bộ phim càng gần với 1 dành cho người lớn. Tương tự trong khoảng  $[-1, 1]$  ta cũng gán mỗi user một giá trị nằm trong khoảng  $[-1, 1]$  đại diện cho tuổi của user. Nếu user càng ít tuổi ta gán giá trị càng gần với -1, ngược lại user càng nhiều tuổi ta gán cho user đó giá trị càng gần với 1.



**Hình 3:** Các dấu checkmark cho ta biết bộ phim mà một user cụ thể đã xem. Từ các dấu checkmark đó ta có thể dễ dàng suy ra được sở thích của người dùng. Khi đó ta có thể biết được những người có độ tuổi trưởng thành thường không thích xem phim dành cho trẻ con và đề xuất phim cho người trưởng thành.

Nguồn:

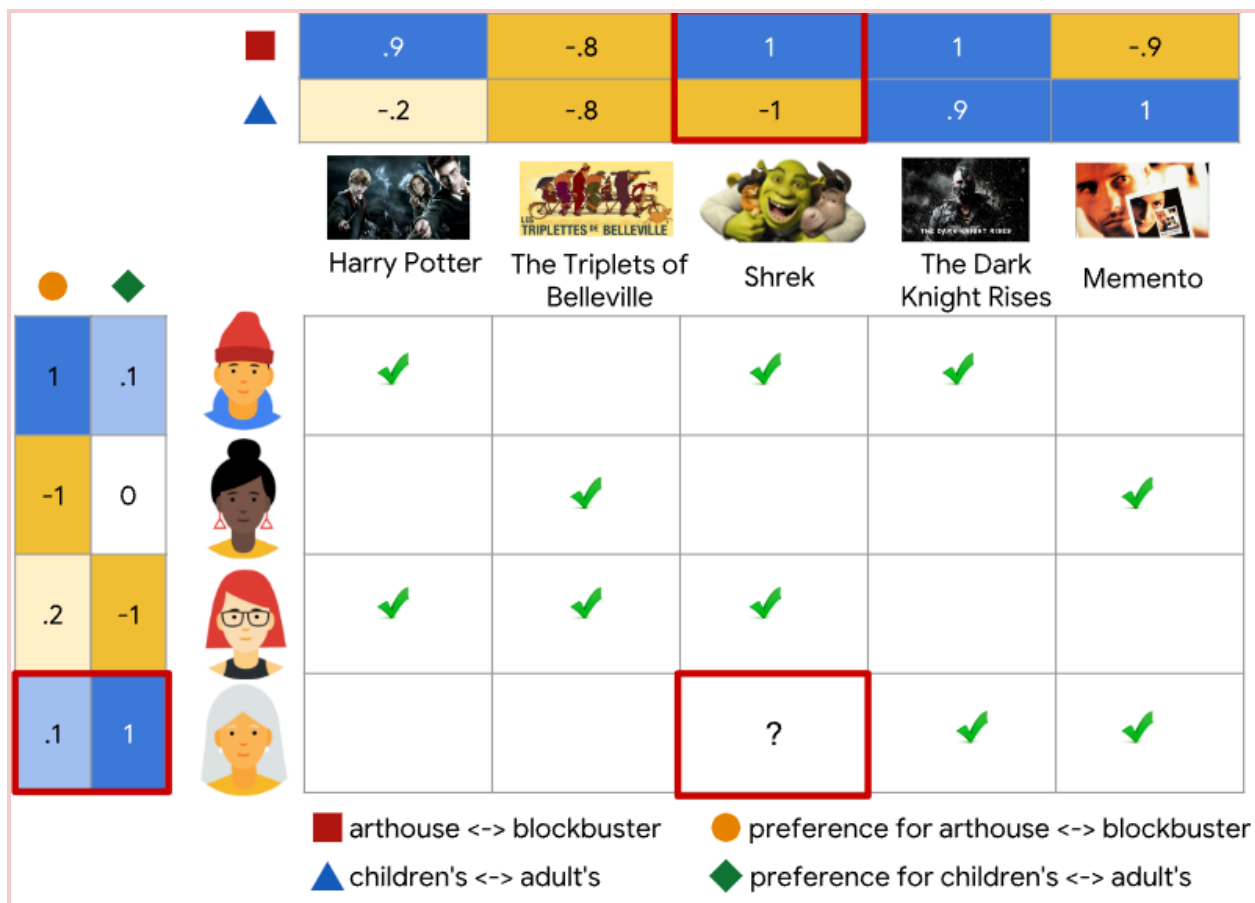
<https://developers.google.com/machine-learning/recommendation/collaborative/basics>

## 2.2.2.2. 2D Embedding

Trong trường hợp một feature không thể đại diện rõ cho từng user và item, ví dụ: ta muốn phân biệt rõ hơn user đầu tiên và user thứ 2 trong Hình 3. Ta có thể thêm 1 feature vào







**Hình 5:** Feedback matrix các dấu checkmark đại diện cho bộ phim mà user đã xem

Nguồn:

<https://developers.google.com/machine-learning/recommendation/collaborative/basics>

Sau khi tạo được feedback matrix, cùng với embedding của mỗi user và item, ta có thể tính toán sự giống nhau giữa các user để từ đó nếu 2 user có mức độ tương đồng cao, hệ thống có thể đề xuất các bộ phim mà người dùng A chưa xem nhưng người dùng B đã xem và người dùng B thích bộ phim.

Để tính toán mức độ tương đồng có rất nhiều cách tính toán, một trong các cách tính toán phổ biến là cosine similarity:

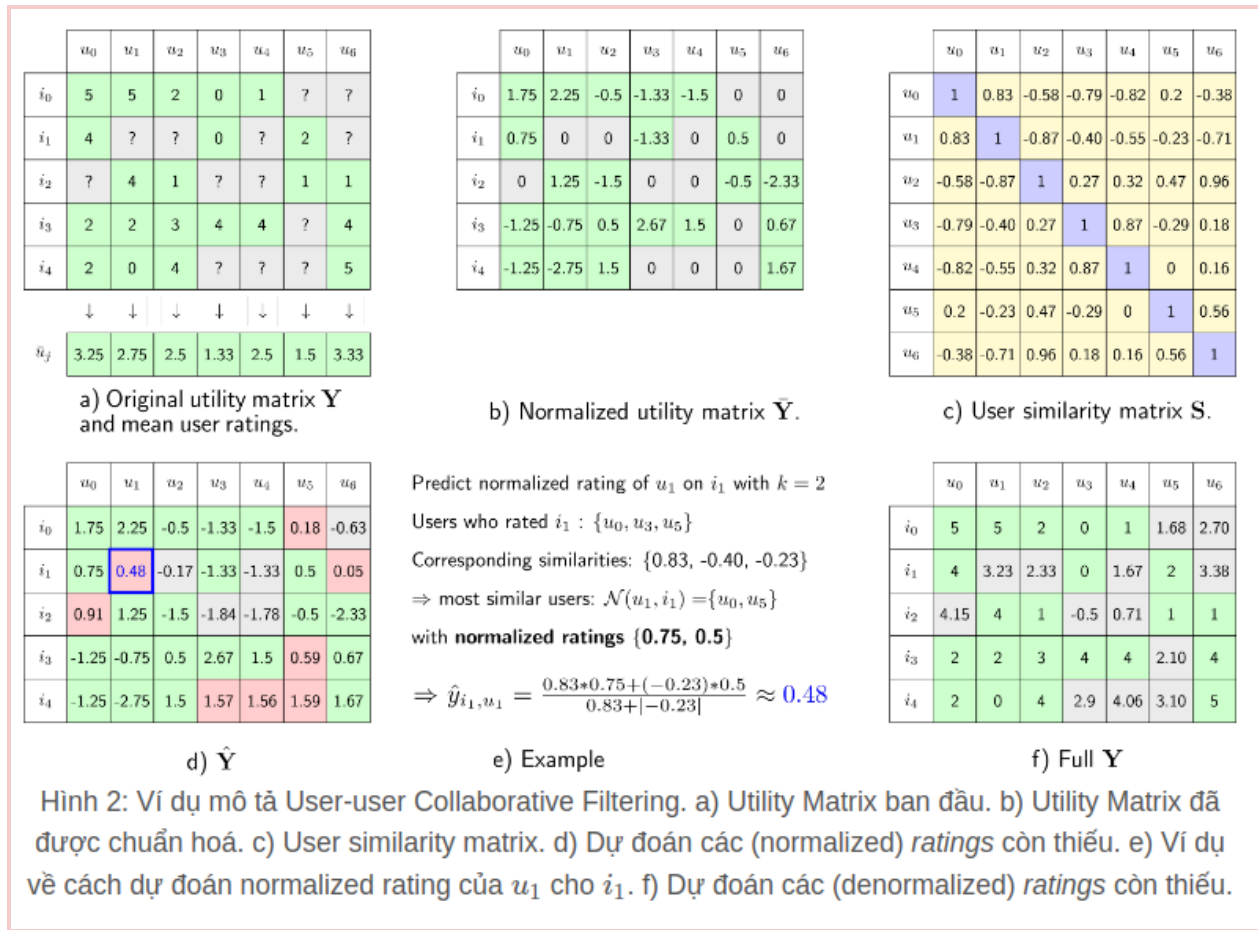
$$\cos(x, y) = \frac{x \cdot y}{||x|| ||y||}$$

Trước khi chuyển sang chi tiết cụ thể trong phương pháp, trong các lần trước ta đã đề cập về vấn đề sự thưa của ma trận rating của cặp ( user, item ). Đây luôn là vấn đề trong bài

toán Recommendation, ở đây để giải quyết vấn đề một cách thiên kiến (bias), chúng tôi đề xuất 2 phương pháp:

1. Thay các giá trị missing rating của cặp (user, item) bằng với giá trị trung bình của các rating mà user đó đã đánh giá. Nhưng khi chọn vấn đề này, có một vài vấn đề khó khăn.
  1. Nếu user đó đã đánh giá rất ít item, ví dụ: user rating 1 item và rating item đó 5 sao, khi đó ta nhận thấy việc giả sử này là không thực tế, và có thể dẫn đến rất nhiều sai sót.
  2. User đó không đánh giá item nào, khi đó ta sẽ không biết điền giá trị gì.
2. Ta sẽ điền vào tất cả cặp giá trị (user, item) missing = 2.5 hoặc 0 tương đương với giá trị trung tính ( neutral ).

### **Example: User-user Collaborative filtering**



**Hình 6:** Mô tả phương pháp User-user Collaborative Filtering.

Nguồn:

<https://machinelearningcoban.com/2017/05/24/collaborativefiltering/#-user-user-collaborative-filtering>

**a)** Là ma trận ban đầu, hàng dưới cùng  $\hat{u}_j$  biểu diễn giá trị trung bình của các cột.

**b)** Normalize matrix với  $u_j = u_j - \hat{u}_j$ .

**c)** Tính toán similarity giữa các user sử dụng **cosine similarity**.

**d)** Là ma trận dự đoán các giá trị missing rating, trong đó các giá trị missing rating được dự đoán bằng cách, chọn ra  $k$  ( trong bài toán này chọn  $k = 2$  ) user có mức độ similarity cao nhất đối với user muốn dự đoán và user đó đã rating item, sau đó ta dự đoán missing rating như trong hình e).

f) Là ma trận sau khi tính toán xong và thay  $u_j = u_j + \hat{u}_j$

### 2.2.2.3 Phương pháp tiếp cận sử dụng MapReduce:

#### Bước 1: Chia dữ liệu thành các cặp khóa - giá trị:

*Mapper:* (key, value)  $\rightarrow$  (user\_id, movie:rating\_score) | map mỗi user với movie:rating\_score tương ứng.

*Reducer:* (key, value)  $\rightarrow$  (user\_id, (movie1: rating\_score, movie2:rating\_score, ..)).

#### Bước 2: Tạo ma trận xuất hiện đồng thời:

*Mapper:* Tạo ma trận xuất hiện đồng thời, ma trận xuất hiện đồng thời ở đây có thể hiểu là đối với mỗi cặp movie xuất hiện đồng thời hay đều được xem bởi user\_id, thì với mỗi cặp movies đó ta tăng giá trị lên 1: (key, value) = (movie1:movie2, <1,1,1,...>)

*Reducer:* Tổng các giá count: (key, value) = (movie1:movie2, sum(<1,1,1...>)).

#### Bước 3: Normalize giá trị quan hệ:

*Mapper:* Định dạng lại output từ các cặp (movie1:movie2, sum()), (movie1:movie3, sum())  $\rightarrow$  (movie1, <movie2:relation, movie3: relation, ...>) (relation = sum).

*Reducer:* Normalize các giá relation thành relation/ sum( relation, .. ) với mỗi key như user 1 trong mapper ở trên sẽ tồn tại nhiều giá trị relation, ta normalize nó bằng cách sum = tổng các giá trị relation đó lại rồi lấy relation\_n = relation/sum  $\Rightarrow$  (movie2, <movie1:realtion\_n, ....>).

#### Bước 4: Tính toán giá trị user rating cho movie qua relation giữa các movies:

Bao gồm 2 quá trình Mapper:

*Mapper1:* Truyền output bước 3 của Reducer trước cho Reducer.

*Mapper2:* Map từ input ban đầu (user\_id, item\_id, rating) thành (movie\_id, user\_id:rating).

*Reducer:* Nhận vào input bao gồm (movie2\_id, <movie1\_id=relation, .. item\_id:rating, ...>)  $\Rightarrow$  (key, value) = (user:movie, relation\*rating).

#### Bước 5: Tính tổng giá trị user rating cho movie:

Mapper : truyền output của Reducer trước trong bước 4 to Reducer.

Reducer: Tính tổng cách giá trị  $relationrating$ ,  $(key, value) \Rightarrow (user:movie, sum())$  [ trong đó  $sum()$  là giá trị rating dự đoán của user dành cho movie ].

**Bước 6: Lọc các giá trị rating của user dành cho movie nhỏ hơn threshold:**

Mapper: Lọc các cặp  $(key, value)$  có value nhỏ hơn threshold  $A$ . Giữ lại các cặp giá trị  $(key, value) = (user:movie, rating)$  lớn hơn hoặc bằng threshold  $A$ .

Reducer: Reduce các cặp  $(key, value) = (user:movie, rating)$  thành  $(user, movie)$  các movie này sẽ được đề xuất cho user.

