

# Hệ thống gợi ý bài báo<sup>1</sup>

\*Phạm Minh Châu, \*\*Lê Thanh Hoàng, \*\*Trần Đình Khang

{ [Chuanpm@gmail.com](mailto:Chuanpm@gmail.com); [huongthanh@gmail.com](mailto:huongthanh@gmail.com); [khangtd-fit@mail.hut.edu.vn](mailto:khangtd-fit@mail.hut.edu.vn) }

\* Học Sĩ Phạm Khắc Thuần, Học Viện

\*\* Viện CNTT & TT - HBKHN

**Tóm tắt.** Hệ thống gợi ý là công cụ hỗ trợ ra quyết định, nhằm mục đích cung cấp cho người sử dụng những gợi ý về thông tin, sản phẩm, và dịch vụ phù hợp nhất với yêu cầu và sở thích riêng của từng người dùng (người dùng) hoặc yêu cầu chung. Có nhiều thuật toán khác nhau xây dựng hệ thống gợi ý, chúng tôi sẽ mô tả các thuật toán, dựa trên nội dung, dựa trên tri thức và kết quả lại. Trong bài viết này chúng tôi sẽ mô tả thuật toán gợi ý bài báo, sử dụng các kỹ thuật khai phá dữ liệu kết hợp với các thuật toán và dựa trên nội dung xây dựng hệ thống gợi ý bài báo phù hợp nhất cho người dùng.

**Từ khóa.** Hệ thống gợi ý, thuật toán, dựa trên nội dung, dựa trên tri thức, kết quả lại, khai phá dữ liệu, thuật toán.

## 1. Giới thiệu và phát biểu bài toán

Ngày nay, người sử dụng các hệ thống thông tin, đặc biệt là các website thương mại, thường khó khăn trong việc tìm kiếm và lựa chọn các thông tin cần thiết và phù hợp với quy định và nhu cầu (ví dụ như việc chọn mua một chiếc ô tô phù hợp, hoặc việc lập kế hoạch cho một chuyến đi du lịch); bởi vì người sử dụng có quá nhiều lựa chọn, nhưng không có thời gian hoặc “tri thức” để đánh giá những lựa chọn này và đưa ra các quyết định.

Hệ thống gợi ý (Recommender Systems) ([8], [12]) là các công cụ hỗ trợ ra quyết định, nhằm mục đích cung cấp cho người sử dụng những gợi ý về thông tin, sản phẩm, và dịch vụ phù hợp nhất với yêu cầu và sở thích riêng của từng người dùng (người dùng) yêu cầu chung. Hiện nay, các hệ thống gợi ý đã trở thành một trong những công cụ hỗ trợ và phổ biến nhất trong các hệ thống thương mại điện tử (ví dụ như Amazon.com, Barnes&Noble.com, Yahoo! news, TripAdvisor.com,...).

Việc xây dựng là những nhà khoa học hoặc những người nghiên cứu thì việc tìm kiếm những bài báo khoa học tham khảo phù hợp với lĩnh vực mình nghiên cứu cũng không phải là việc dễ dàng. Những phân tích trên nhóm sẽ đưa ra một hệ thống gợi ý bài báo đáp ứng nhu cầu tìm kiếm tài liệu tham khảo bằng tin tức cho các nhà khoa học, các nhà nghiên cứu tại Việt Nam.

Bài toán gợi ý bài báo có thể phát biểu như sau. Hệ thống bao gồm người dùng, và những bài báo. Khi người dùng có yêu cầu tìm kiếm các bài báo bằng việc nhập thông tin cần tìm kiếm (các từ khóa) sau đó hệ thống sẽ đưa ra danh sách các bài báo mà theo đó đánh giá các bài báo này phù hợp nhất với người dùng (có mức độ ưu tiên cao với người dùng).

Phát biểu bài toán bằng công thức toán học như sau:

Giả sử

$U$  là tập người dùng trong hệ thống.

$I$  là tập bài báo trong hệ thống.

$u$  là một người dùng trong  $U$

$i$  là một bài báo trong  $I$

$r_{ui}$  là đánh giá của người dùng  $u$  cho bài báo  $i$ .  $r_{ui}$  nhận giá trị trong tập  $X \subset \mathbb{R}$ . Trong đó, tập  $X$  là các số nguyên  $\{1, 2, 3, 4, 5\}$  tương ứng với các mức đánh giá Ghét, Không thích, , thích, rất thích. .

Bài báo có biểu diễn nội dung có cấu trúc như sau:

Item = (PaperID, Title, Date, Magazine, Authors, Keyword1, Keyword2, Domain, References, Content)

✓ PaperID: ID của bài báo

✓ Title: tiêu đề bài báo, có biểu diễn bằng một vector trong số các từ khóa.

<sup>1</sup> Bài báo hoàn thành với tài trợ của Quỹ Phát triển Khoa học và Công nghệ Quốc gia (Nafosted), mã số tài 102.01.30.09.

- ✓ Date: ngày đăng của bài báo.
- ✓ Magazine: tạp chí đăng.
- ✓ Authors: Tập thể tác giả, biểu diễn bằng vector Id của các tác giả.
- ✓ Keyword1: tập các từ khóa xuất hiện trong Keyword của bài báo, biểu diễn bằng vector trọng số của các từ khóa.
- ✓ Keyword2: tập các từ khóa đi kèm cho toàn bộ bài báo
- ✓ Domain: lĩnh vực của bài báo.
- ✓ References: tập các bài báo tham khảo, biểu diễn bằng vector Id của các bài báo.
- ✓ Content: nội dung của bài báo.

Thông tin người dùng biểu diễn như sau:

User = (UserID, Name, sex, Domains, Keywords, Id\_author)

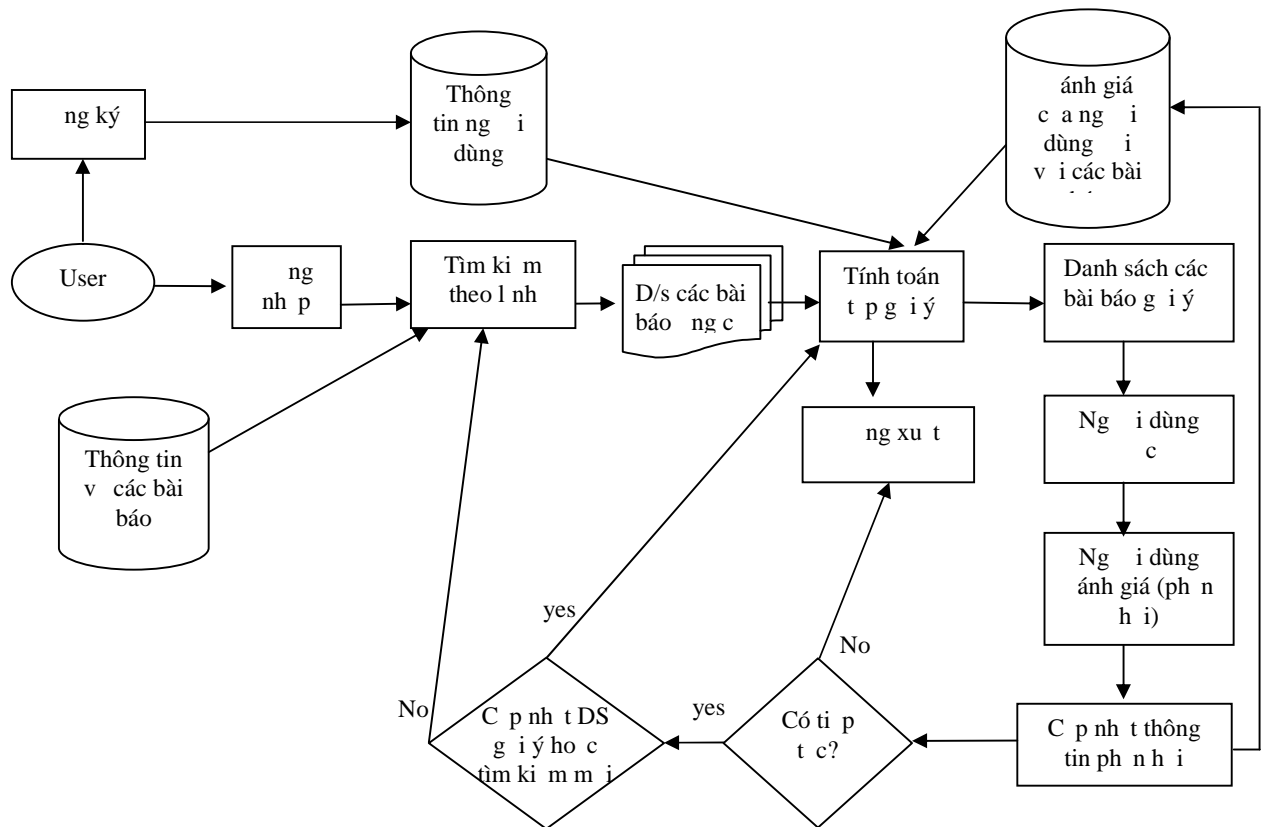
- ✓ Name: tên của người dùng
- ✓ Sex: giới tính
- ✓ Domains: nghề nghiệp của người dùng
- ✓ Keywords: những từ khóa yêu thích (tên sự kiện tìm kiếm)
- ✓ Id\_author: tác giả yêu thích (tên sự kiện bài báo của tác giả)

Bảng đánh giá người dùng về các bài báo

- ✓ PaperID
- ✓ UserID
- ✓ Rate

## 2. Mô hình hệ thống gợi ý bài báo

Mô hình tổng quát của hệ thống gợi ý trình bày trong hình 1 dưới đây.



Hình 1. Mô hình hệ thống gợi ý bài báo

Khi người dùng sử dụng hệ thống, nhu cầu người ký thì hệ thống yêu cầu người ký và cung cấp một số thông tin cá nhân cần thiết. Khi người nhập vào hệ thống; người dùng nhập các từ khóa để tìm kiếm; hệ thống dựa trên thông tin về người dùng để đề xuất các từ khóa mà người dùng nhập vào để tìm kiếm trong kho dữ liệu về các bài báo (tìm trong keyword2) dựa trên các bài báo người dùng, sau đó hệ thống sẽ tính toán top k các bài báo mà phù hợp nhất với người dùng dựa trên các từ khóa và các tác giả yêu thích; và sẽ đề xuất các bài báo người dùng và những bài báo mà người dùng đã thích để đánh giá kết quả và sẽ đánh giá các bài báo người dùng có cùng sở thích với người dùng dựa trên các bài báo người dùng; sau đó hệ thống sẽ hiển thị danh sách các bài báo phù hợp nhất cho người dùng, người dùng click và nhìn liệu hệ thống các bài báo đó; hệ thống sẽ tiến hành cập nhật thông tin người dùng; người dùng nhập từ khóa thì sẽ có hai lựa chọn hoặc là cập nhật danh sách gợi ý hoặc là tiến hành tìm kiếm mới và hệ thống sẽ tiếp tục theo nhu cầu của người dùng.

Dựa trên các dữ liệu thu thập là đánh giá của người dùng về các bài báo trong hệ thống (explicit data) kết hợp với thông tin cá nhân của người dùng như là lịch sử xem, danh sách các bài báo đã đọc, áp dụng các mô hình tính toán, để đoán đánh giá của người dùng xác định vị trí bài báo thích hợp. Có thể trong hệ thống này chúng tôi áp dụng kỹ thuật cộng tác (Collaborative filter), dựa trên đánh giá của người dùng về các bài báo; và kỹ thuật dựa trên nội dung (Content based) kết hợp với đánh giá người dùng về các bài báo, mà phù hợp để cập nhật bài báo với người dùng dựa trên danh sách các bài báo phù hợp cho người dùng trong hệ thống.

### Một số tình huống của hệ thống:

**Tình huống thực tế:** Khi có một người sử dụng mới người này chỉ có thông tin cá nhân cần thu thập khi người ký; chưa có bất kỳ đánh giá nào về các bài báo để nên khó khăn cho hệ thống gợi ý tính toán những bài báo phù hợp nhất với người dùng.

*Gợi ý pháp:*

Dựa trên những từ khóa tìm kiếm của người dùng, và những từ khóa yêu thích tính mức độ phù hợp của bài báo với người dùng và dựa trên những bài báo thích thú cho người dùng đó; khi người dùng click bài báo và có phản hồi thì hệ thống sẽ phân tích dựa trên gợi ý mới.

**Tình huống thực tế hai:** Khi có một bài báo mới mà chưa có người dùng nào đánh giá hoặc có ít người dùng đánh giá, lúc đó hệ thống sẽ xử lý những thông tin gợi ý cho người dùng tiếp theo.

*Gợi ý pháp:*

Sử dụng thông tin về bài báo như là lịch sử; thông tin về tác giả bài báo; sẽ đề xuất các bài báo dựa trên gợi ý cho người dùng trong hệ thống. Có thể, sử dụng thông tin về tác giả bài báo (tác giả này có thể là một người dùng trong hệ thống) chúng ta có thể tìm ra những người dùng có sở thích tương tự như tác giả bài báo và để đoán đánh giá của người dùng sử dụng đó về bài báo; hoặc chúng ta có thể dựa vào sẽ đề xuất người dùng các bài báo đã đánh giá bình thường sử dụng nào đó về bài báo mới đó có thể tính toán xem bài báo này có phù hợp nhất với người dùng tiếp theo, ...

**Tình huống thực tế ba:** Với người dùng đã có trong hệ thống, chúng ta có thể áp dụng kỹ thuật cộng tác, kỹ thuật dựa trên nội dung kết hợp các kỹ thuật khai phá dữ liệu dựa trên gợi ý phù hợp nhất với người dùng.

## 3. Các kỹ thuật sử dụng trong quá trình gợi ý

### 3.1 Quá trình tìm kiếm lý tưởng

Trước khi đi vào các thuật toán top k gợi ý các bài báo thì hệ thống cần phải tìm kiếm lý tưởng các dữ liệu đầu vào; đây chính là top các bài báo. Với vị trí bài báo lưu trữ để định lượng và phân bố, và với việc thi triển các bài báo đã phân theo lịch sử xem của người dùng.

Với vị trí bài báo chúng ta sẽ tiến hành phân tách các từ theo phương pháp Maximum Matching. Đây là phương pháp tách từ ngắn và dài nhất. Phương pháp này sẽ định lượng và phân bố làm các phân tách các từ.

Sau khi đã phân tách các từ chúng ta sẽ tiến hành lựa chọn các từ cho từng văn bản. Lựa chọn các từ là vì các từ không mang thông tin khi văn bản nhằm nâng cao hiệu quả phân loại và giảm số

ph c t p tính toán. Có nhi u ph ng pháp l a ch n c tr ng, tuy nhiên trong khuôn kh bài vi t này chúng tôi s d ng ph ng pháp Information Gain. Các c tr ng sau khi ã c trích ch n c l u trong Keyword2 c a bài báo.

Information Gain xác nh l ng thông tin c a m t t c s d ng cho vi c phân lo i d a trên s có m t hay không có m t c a t ó trong m t v n b n.

Bài toán l a ch n c tr ng c a v n b n c phát bi u nh sau: cho  $c_1, \dots, c_k$  là t p h p g m k lo i (hay k l nh v c).  $P(c_j)$  là xác su t tìm c m t v n b n có trong t p d li u và v n b n ó thu c lo i  $c_j$ .  $P(w)$  là xác su t xu t hi n c a t w.  $P(w, c_j)$  là xác su t tìm c m t v n b n thu c lo i  $c_j$  và v n b n ó ch a t w.  $P(\bar{w}, c_j)$  là xác su t tìm c m t v n b n thu c lo i  $c_j$  và v n b n ó không ch a t w.

Information Gain c a m t t w c tính nh sau:

$$IG(w, c_i) = \sum_{c \in \{c_i, \bar{c}_i\}} \sum_{w \in \{w_k, \bar{w}_k\}} P(w, c) \cdot \log \frac{P(w, c)}{P(w) \cdot P(c)} \quad (1)$$

Theo ph ng pháp này ng i ta d a vào các hàm tính toán  $IG(w, c_i)$  c a t w i v i phân lo i  $c_i$  quy t nh xem có nên l a ch n t w làm c tr ng c a tài li u hay không. Ta c n tính toán IG cho m i t trong t p d li u hu n luy n và nh ng t có IG nh h n m t ng ng cho tr c s b lo i b .

bi u di n các bài báo chúng ta s d ng vector trong không gian K chi u; chúng ta có th s d ng mô hình t n su t bi u di n s d ng ph ng pháp **TF x IDF**:

### Ph ng pháp TF x IDF

V i bài báo  $d_j$  c bi u di n b i m t vector c tr ng  $d_j(w_{1j}, w_{2j}, \dots, w_{kj})$

Trong ó  $w_{ij}$  là tr ng s c a t khóa  $t_i$  trong bài báo  $d_j$

Ph ng pháp này là t ng h p c a hai ph ng pháp TF và IDF ([11]), giá tr c a ma tr n tr ng s c tính nh sau:

$$w_{ij} = \begin{cases} \left[ 1 + \log\left(\frac{f_{ij}}{\max_z f_{zj}}\right) \right] \log\left(\frac{m}{h_i}\right) & \text{if } h_i \geq 1 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Trong ó,  $f_{ij}$  là s l n xu t hi n c a t khóa  $t_i$  trong bài báo  $d_j$  và  $\max_z f_{zj}$  là s l n xu t hi n l n nh t c a m t t khóa trong tài li u  $d_j$ ,  $m$  là s l ng bài báo và  $h_i$  là s bài báo mà t khóa  $t_i$  xu t hi n.

Ph ng pháp này k t h p c u i m c a c hai ph ng pháp TF và IDF. Tr ng s  $w_{ij}$  c tính b ng t n su t xu t hi n c a t khóa  $t_i$  trong bài báo  $d_j$  và  $h_i$  m c a t khóa  $t_i$  trong toàn b c s d li u.

Vi c gán tr ng s cho các t khóa c tr ng trong Keyword2 s xem xét n các t khóa xu t hi n trong tiêu bài báo và trong Keyword1 (do ng i dùng t nh p), v i các t khóa này thì c v i tr ng s có m c quan tr ng h n so v i các t khóa c tr ng khác.

## 3.2. M t s thu t toán liên quan

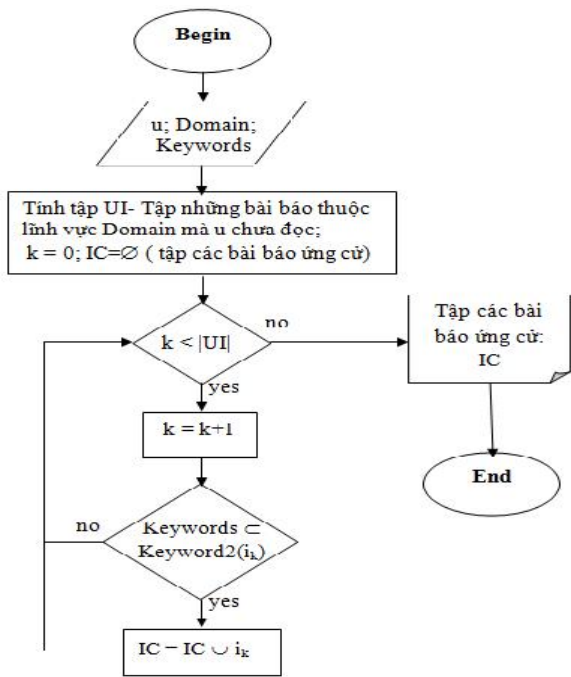
### 3.2.1 Thu t toán a ra danh sách các bài báo ng c

u vào: danh sách các t khóa, các bài báo trong h th ng v i keyword2 t ng ng

u ra: danh sách các bài báo ng c

Quá trình x lý:

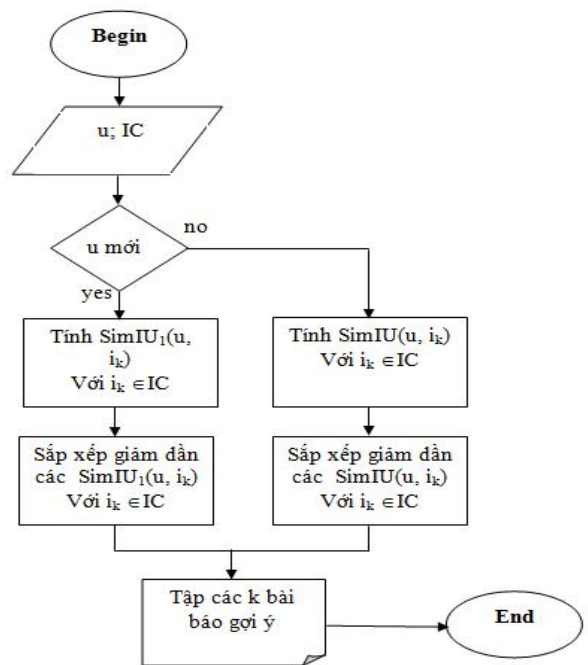
- Tìm ki m các bài báo có ch a các t khóa c n tìm trong Keyword2
- N u không tìm th y yêu c u ng i dùng nh p l i t khóa



Hình 2. Thuật toán tìm ra danh sách các bài báo ứng cử

**3.2.2 Thuật toán tìm ra danh sách các bài báo gợi ý**

u vào: Danh sách các bài báo ứng cử  
 u ra: Danh sách k bài báo phù hợp nhất gợi ý cho người dùng  
 Quá trình xử lý:



Hình 3. Thuật toán tìm ra danh sách các bài báo gợi ý

Trường hợp thứ nhất: nếu người dùng mới sử dụng thì chúng ta sử dụng công thức (6) tính toán mức phù hợp của các bài báo ứng cử với người dùng dựa trên những từ khóa yêu thích để tìm ra k bài báo phù hợp nhất với người dùng.

Trường hợp thứ 2: khi người dùng đã sử dụng thì chúng ta sẽ sử dụng công thức (9) tính mức phù hợp của bài báo ứng cử với người dùng để tìm ra k bài

báo phù hợp nhất với ngữ cảnh.; trong công thức này vì cần xác định các hệ số phù hợp để quan trọng; chúng tôi sử dụng hành xác định giá trị của các hệ số phù hợp qua thực nghiệm.

Vấn đề cần khắc phục bao nhiêu là thích hợp chúng tôi sử dụng hành thực nghiệm xác định.

### 3.3 Kỹ thuật đánh giá tương tự giữa hai bài báo

Vì cần đánh giá tương tự của hai bài báo để quan trọng; trong bài viết này chúng tôi sử dụng đánh giá tương tự giữa hai bài báo kết hợp hai cách sau:

Cách 1: dựa trên các từ khóa trọng và trọng số của nó trong hai bài báo

Cách 2: dựa trên số đánh giá của cùng một từ ngữ dùng trong hai bài báo có

Vì cách thứ nhất; chúng tôi dựa trên hai tiêu chí:

- ✓ Số từ khóa giống nhau của hai bài báo (mỗi bài báo sẽ có biểu diễn bằng một tập các từ khóa trọng và trọng số từ ngữ của từ khóa)
  - ✓ Mức quan trọng của các từ khóa trong hai bài báo (thể hiện bằng trọng số của từ khóa trong bài báo)
- đánh giá tương tự của hai bài báo dựa trên các từ khóa trọng, chúng tôi sử dụng công thức sau:

$$Sim_1(i, j) = \frac{m}{K} \left( 1 - \frac{\sum_{k=1}^m |w_{ki} - w_{kj}|}{\sum_{k=1}^m (w_{ki} + w_{kj})} \right) \quad (3)$$

Trong đó m là tổng số từ khóa giống nhau của hai bài báo; K là tổng số từ khóa của bài báo;  $w_{ki}$  là trọng số của từ khóa thứ k trong bài báo i.

Ý tưởng cơ bản này là, trong hai bài báo khác nhau nếu có cùng nhiều từ khóa giống nhau thì tương đồng càng cao, ngoài ra đánh giá chính xác hơn chúng ta quan tâm đến trọng số của mỗi từ khóa giống nhau trong hai bài báo, nếu trọng số của các từ khóa càng gần nhau thì mức tương đồng giữa hai bài báo càng cao.

Vì cách thứ 2; chúng tôi sử dụng tương tự PC (Pearson Correlation) [1].

$$Sim_2(i, j) = \frac{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)(r_{uj} - \bar{r}_j)}{\sqrt{\sum_{u \in U_{ij}} (r_{ui} - \bar{r}_i)^2 \sum_{u \in U_{ij}} (r_{uj} - \bar{r}_j)^2}} \quad (4)$$

Trong đó:

$U_{ij}$ : tập từ ngữ dùng cùng đánh giá hai bài báo i, j

$r_{ui}$  và  $r_{uj}$ : đánh giá của ngữ cảnh dùng cho mỗi bài báo i, j

$\bar{r}_i, \bar{r}_j$ : đánh giá trung bình trong bài báo i, j xét trên tập từ ngữ đánh giá của ngữ cảnh dùng trong tập  $U_{ij}$ .

Công thức tính tương tự tổng quát cho hai bài báo sẽ như sau:

$$Sim(i, j) = w_1 * Sim_1(i, j) + w_2 * Sim_2(i, j) \quad (5)$$

Trong đó  $w_1, w_2$  là các hệ số tương tự cần xác định thông qua thực nghiệm, thỏa mãn điều kiện  $w_1, w_2 \geq 0$ ; và  $w_1 + w_2 = 1$ .

Công thức 5 sẽ sử dụng tính mức phù hợp giữa bài báo với ngữ cảnh trình bày mục 3.5.

### 3.4 Kỹ thuật đánh giá mức tương tự giữa hai ngữ cảnh dùng

đánh giá mức tương tự giữa hai ngữ cảnh dùng chúng tôi sử dụng tương tự PC, tương tự như công thức 4.

$$SimU(u, v) = \frac{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in I_{uv}} (r_{ui} - \bar{r}_u)^2 \sum_{i \in I_{uv}} (r_{vi} - \bar{r}_v)^2}} \quad (6)$$

Trong đó:

$I_{uv}$ : tập tất cả các bài báo có ảnh hưởng của hai người dùng  $u, v$

$r_{ui}$  và  $r_{vi}$ : ảnh hưởng của người dùng  $u$  cho bài báo  $i$

$\bar{r}_u, \bar{r}_v$ : ảnh hưởng trung bình của người dùng  $u, v$  trong tập các bài báo  $I_{uv}$

### 3.5 Kỹ thuật ảnh hưởng của phù hợp của một bài báo với một người dùng

Việc ảnh hưởng của phù hợp của một bài báo trong danh sách đề xuất người dùng sẽ giúp cho quá trình lọc các bài báo gợi ý cho người dùng cá nhân trở nên thuận tiện hơn. Chúng tôi sử dụng phương pháp ảnh hưởng của phù hợp giữa một bài báo với người dùng bằng việc kết hợp ba cách sau:

**Cách 1:** dựa trên mức quan trọng của các từ khóa mà người dùng yêu thích trong bài báo; nếu bài báo nào càng chứa nhiều từ khóa yêu thích và các từ khóa đó có mức quan trọng càng cao với bài báo đó thì bài báo đó sẽ có mức phù hợp càng cao với người dùng.

Công thức tính phù hợp theo cách này có như sau:

$$SimU_1(u, i) = \frac{n}{m} \sum_{k=1}^m w_{ki} \quad (7)$$

Trong đó  $m$  là tổng số từ khóa người dùng yêu thích,  $n$  là số từ khóa mà người dùng yêu thích xuất hiện trong bài báo  $i$ ;  $w_{ki}$  là trọng số của từ khóa  $k$  xuất hiện trong bài báo  $i$ .

Công thức này có thể dùng trong trường hợp người dùng mới đăng nhập hệ thống và ảnh hưởng của bài báo.

**Cách 2:** dựa trên ý tưởng của kỹ thuật lọc đề xuất.

Đây là phương pháp quen thuộc, và sử dụng nguyên tắc Hệ thống đề xuất đề xuất (Collaborative Recommender System) tập hợp thông tin về ảnh hưởng, hoạt động của người dùng; xác định mức độ ảnh hưởng của người dùng dựa trên sở thích, tổng cộng các gợi ý kết quả so sánh trên.

Chúng ta sử dụng công thức (6) tính toán mức phù hợp giữa người dùng người dùng đó dựa trên người dùng (ký hiệu UN) có mức phù hợp cao nhất với người dùng  $u$ .

Sau đó chúng ta tiến hành xét tập người dùng UN có ảnh hưởng của bài báo  $i$  bằng xét, xác định ảnh hưởng trung bình của bài báo  $i$  xét trên tập tất cả những ảnh hưởng của người dùng trong tập UN dựa trên phù hợp.

Công thức tính phù hợp có như sau:

$$SimU_2(u, i) = \frac{\sum_{un \in UN} r_{un,i}}{5 * |UN|} \quad (8)$$

Trong đó  $r_{un,i}$  là ảnh hưởng của người dùng  $un$  với bài báo  $i$ ; đây chúng ta sử dụng thang ảnh hưởng có 5 mức do vậy chuẩn hóa mức phù hợp (nhân giá trị từ 0 đến 1) của bài báo  $i$  với người dùng  $u$  chúng ta sẽ chia cho 5.

**Cách 3:** dựa trên ý tưởng của phương pháp lọc theo nội dung; với cách này chúng ta sẽ ảnh hưởng của phù hợp của bài báo  $i$  với người dùng  $u$ , bằng cách ảnh hưởng của tổng trung bình giữa người dùng bài báo mà người dùng  $u$  đã từng đọc trong quá khứ với bài báo  $i$ .

Giống với IP là tập người dùng bài báo mà người dùng  $u$  đã có ảnh hưởng trong quá khứ; khi đó chúng ta sử dụng công thức (5) tính tổng giá trị bài báo  $i$  với tập người dùng bài báo trong IP và tính trung bình tổng của bài báo  $i$  với các bài báo đó.

Công thức tính phù hợp trong trường hợp này có xu hướng sau:

$$SimIU_3(u, i) = \frac{\sum_{j \in IP} Sim(i, j)}{|IP|} \quad (9)$$

Công thức tổng quát mà chúng tôi xu hướng tính toán mô phù hợp của bài báo  $i$  vì vì  $i$  dùng  $u$  như sau:

$$SimIU(u, i) = l_1 * SimIU_1(u, i) + l_2 * SimIU_2(u, i) + l_3 * SimIU_3(u, i) \quad (10)$$

Trong đó  $l_1, l_2, l_3$  là các hệ số phù hợp tính toán trong quá trình nghiên cứu;  $l_1, l_2, l_3 \geq 0$  và  $l_1 + l_2 + l_3 = 1$ .

### 3.6. Cài đặt ứng dụng và đánh giá kết quả thực nghiệm

Chúng tôi đã tiến hành cài đặt thí nghiệm với từng số mặt từ bài báo thu được 81 như v: Công nghệ mạng và truy cập thông tin, Trí tuệ nhân tạo, Xử lý ngôn ngữ tự nhiên và tiếng nói, Công nghệ phần mềm, Cơ sở dữ liệu và hệ thống thông tin, Cơ sở toán học của CNTT, Khoa học hệ thống và quản lý, Các phương pháp tính toán mô.

Các bài báo trong hệ thống đã được đánh giá bởi 50 người bao gồm các chuyên gia và các nhà nghiên cứu trong lĩnh vực Công nghệ Thông tin và toán tin, bao gồm tất cả các kết quả quan trọng, lúc ban đầu là 20 người có ý kiến đưa vào nội dung khóa tìm kiếm, đưa vào lĩnh vực nghiên cứu trên từng bài báo 45, chính xác 10%; sau khi hệ thống có dữ liệu đánh giá của người dùng thì kết quả thí nghiệm với số lượng 40 người dùng và 70 bài báo thì chính xác tăng 19% và khi thí nghiệm với 50 người dùng và 100 bài báo thì chính xác của quá trình tăng 25%.



Hình 4. Giao diện chính của hệ thống gợi ý bài báo



#### 4. Kết luận và hướng phát triển

Trong bài báo chúng tôi đã xuất k thu t ánh giá t ng ng c a hai bài báo đ a trên nh ng t khóa c tr ng và đ a trên s ánh giá c a cùng m t t p ng i dùng v i hai bài báo ó; trong k thu t này c n ph i xác nh các h s t ng t  $w_1$  và  $w_2$  làm sao cho t hi u qu cao nh t. Ngoài ra bài vi t còn xu t thêm k thu t ánh giá phù h p gi a m t bài báo i v i m t ng i dùng s đ ng hai ý t ng ch o c a ph ng pháp l c c ng tác và đ a trên n i dung. Trong k thu t này vi c xác nh các h s phù h p c ng r t quan tr ng. Trong h th ng c a chúng tôi, các h s này c xác nh qua th c nghi m.

H th ng g i ý bài báo mà chúng tôi xu t b c u ã c cài t th nghi m, tuy nhiên ch a c ki m ch ng b i m t t p l n đ li u trong th c t , do v y chính xác ch a cao. Chúng tôi s ph i làm vi c nhi u h n n a hoàn thi n mô hình và ã ra th nghi m v i kh i l ng l n các bài báo và ng i dùng.

## 5. Tài li u tham kh o

1. Francesco Ricci, Lior Rokach, Bracha Shapira, Paul B.Kantor,: *Recommender Systems Handbook* , Springer, 13 Dec 2010.
2. Gediminas Adomavicius and Alexander Tuzhilin: *Toward the Next Generation of Recommender Systems state of the art 2005*, IEEE, 2005.
3. Stefan B. Beckers,: *Recommender system*, Duisburg Publisher ,2006.
4. Matthew R. McLaughlin and Jonathan L. Herlocker , *A Collaborative Filtering Algorithm and Evaluation Metric* , Oregon Univer Publisher, 2004.
5. Daniel Billsus and Michael J. Pazzani , *Learning Collaborative Information Filters*.
6. Robin Burke , Hybrid Recommender Systems.
7. Thomas Tran and Robin Cohen, *Hybrid Recommender Systems for Electronic Commerce*, AAAI, 2000.
8. Resnick, P., & Varian, H. R. (1997), *Recommender systems. Communications of the ACM* , 40 (3), 56-58.
9. Two Crows, *Introduction to Data Mining and Knowledge Discovery*, Third Edition, <http://www.twocrows.com/booklet.htm>.
10. Fabrizio Sebastiani, *Machine Learning in Automated Text Categorization*, ACM Computing Surveys, Vol. 34, No. 1, March 2002.
11. Salton, G., & McGill, M. J. (1986), *Introduction to Modern Information Retrieval*, McGraw-Hill, Inc.
12. Billsus, D., & Pazzani, M. J. (1999), *A personal news agent that talks, learns and explains. AGENTS '99: Proceedings of the third annual conference on Autonomous Agents*, (pp. 268-275). ACM.
13. Adomavicius, G., & Tuzhilin, A. (2005), *Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING* , 17, 734-74.

Ng

N  
á