

DS300.N11 - Hệ Khuyến Nghị

Bài tập khuyến nghị dựa trên nội dung

Nguyễn Văn Kiệt, Huỳnh Văn Tín
Sinh viên: Phạm Đức Thế - 19522253

Thứ 3, ngày 27 tháng 09 năm 2022

Bài tập

Khuyến nghị dựa trên nội dung

Bộ dữ liệu phim

ID phim	Nội dung chính
M1	Cuộc chiến đấu chống lại linh hồn ác quỷ
M2	Cuộc đột nhập của Thom vào một ngân hàng tại Hồng Kông
M3	Những rắc rối hài hước đời thường của cô gái 30
M4	Siêu anh hùng cùng chống lại mối nguy hiểm

Lịch sử xem phim của User1

Lịch sử xem phim	ID phim	Nội dung chính
User1	M5	Tình yêu hài hước của cô gái biên kịch

Tiền xử lý dữ liệu

- Tách từ

ID phim	Nội dung chính
M1	Cuộc chiến_đấu chống lại linh_hồn ác_quỷ
M2	Cuộc đột_nhập của Thom vào một ngân_hàng tại Hồng_Kông
M3	Những rắc_rối hài_hước đời_thường của cô_gái 30
M4	Siêu anh_hùng cùng chống lại mối nguy_hiểm

Lịch sử xem phim	ID phim	Nội dung chính
User1	M5	Tình_yêu hài_hước của cô_gái biên_kịch

- Chuyển nội dung về từ thường

ID phim	Nội dung chính
M1	cuộc chiến_đấu chống lại linh_hồn ác_quỷ
M2	cuộc đột_nhập của Thom vào một ngân_hàng tại Hồng_Kông
M3	những rắc_rối hài_hước đời_thường của cô_gái 30
M4	siêu anh_hùng cùng chống lại mối nguy_hiểm

Lịch sử xem phim	ID phim	Nội dung chính
User1	M5	tinhh_yeuhài_huoc của cô_gáibiên_kich

- Loại bỏ từ trùng

Từ điển = {cuộc, chiến_đấu, chống, lại, linh_hồn, ác_quỷ, đột_nhập, của, Thom, vào, một, ngân_hàng, tại, Hồng_Kông, những, rắc_rối, hài_hước, đời_thường, cô_gái, 30, siêu, anh_hùng, cùng, mối, nguy_hiểm}

1. Chuyển tất cả các bộ phim về vector Boolean và sử dụng cosine để tìm ra top 1 bộ phim khuyến nghị cho user1.

- Vector Boolean:

- M1 = {1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0}
- M2 = {1, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0}
- M3 = {0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0}
- M4 = {0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1}
- M5 = {0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0}

- Tính toán sử dụng cosine:

- Công thức:

$$\cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| \cdot |\vec{b}|} = \frac{a_1 \cdot b_1 + a_2 \cdot b_2}{\sqrt{a_1^2 + a_2^2} \cdot \sqrt{b_1^2 + b_2^2}}$$

- Áp dụng:

$$\begin{aligned} \cos(\vec{M5}, \vec{M1}) &= \frac{\vec{M5} \cdot \vec{M1}}{|\vec{M5}| \cdot |\vec{M1}|} = \frac{0}{\sqrt{3} \cdot \sqrt{6}} = 0 \\ \cos(\vec{M5}, \vec{M2}) &= \frac{\vec{M5} \cdot \vec{M2}}{|\vec{M5}| \cdot |\vec{M2}|} = \frac{1}{\sqrt{3} \cdot \sqrt{9}} = \frac{\sqrt{3}}{9} \\ \cos(\vec{M5}, \vec{M3}) &= \frac{\vec{M5} \cdot \vec{M3}}{|\vec{M5}| \cdot |\vec{M3}|} = \frac{3}{\sqrt{3} \cdot \sqrt{7}} = \frac{\sqrt{21}}{7} \\ \cos(\vec{M5}, \vec{M4}) &= \frac{\vec{M5} \cdot \vec{M4}}{|\vec{M5}| \cdot |\vec{M4}|} = \frac{0}{\sqrt{3} \cdot \sqrt{7}} = 0 \end{aligned}$$

- Kết luận: Top 1 bộ phim khuyến nghị cho user1 là bộ phim có ID là M3.

2. Chuyển tất cả các bộ phim về vector TF-IDF và sử dụng cosine để tìm ra top 1 bộ phim khuyến nghị user1.

Biểu diễn nội dung - content(p), TF-IDF

Đưa vào một từ t và một văn bản d

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \cdot \text{IDF}(t)$$

□ IF: Term Frequency

- Tần suất xuất hiện của từ t trong văn bản d
- Những từ quan trọng xuất hiện thường xuyên hơn
- $\text{TF}(t, d) = \text{count of } t \text{ in } d / \text{number of words in } d$ (hoặc có thể không chia)

□ IDF: Inverse Document Frequency

- Giảm trọng số cho những từ xuất hiện trong hầu hết các tài liệu.

$$\text{IDF}(t) = \log\left(\frac{N}{df(t) + 1}\right)$$

- N: Tổng số tài liệu
- $df(t)$: Số tài liệu mà từ t xuất hiện trong N tài liệu

- Vector TF-IDF:

- M1 = {0.018, 0.043, 0.018, 0.018, 0.043, 0.043, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0}
- M2 = {0.014, 0, 0, 0, 0, 0, 0.033, 0.014, 0.033, 0.033, 0.033, 0.033, 0.033, 0.033, 0, 0, 0, 0, 0, 0, 0, 0, 0}
- M3 = {0, 0, 0, 0, 0, 0, 0.018, 0, 0, 0, 0, 0, 0.043, 0.043, 0.043, 0.043, 0.043, 0, 0, 0, 0, 0, 0}
- M4 = {0, 0, 0.018, 0.018, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0.043, 0.043, 0.043, 0.043, 0.043, 0}
- M5 = {0, 0, 0, 0, 0, 0, 0.032, 0, 0, 0, 0, 0, 0, 0, 0, 0.074, 0, 0.074, 0, 0, 0, 0, 0, 0}

- Tính toán sử dụng cosine:

– Công thức:

$$\cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| \cdot |\vec{b}|} = \frac{a_1 \cdot b_1 + a_2 \cdot b_2}{\sqrt{a_1^2 + a_2^2} \cdot \sqrt{b_1^2 + b_2^2}}$$

– Áp dụng:

$$\cos(\vec{M5}, \vec{M1}) = \frac{\vec{M5} \cdot \vec{M1}}{|\vec{M5}| \cdot |\vec{M1}|} = 0$$

$$\cos(\vec{M5}, \vec{M2}) = \frac{\vec{M5} \cdot \vec{M2}}{|\vec{M5}| \cdot |\vec{M2}|} = 0.046$$

$$\cos(\vec{M5}, \vec{M3}) = \frac{\vec{M5} \cdot \vec{M3}}{|\vec{M5}| \cdot |\vec{M3}|} = 0.593$$

$$\cos(\vec{M5}, \vec{M4}) = \frac{\vec{M5} \cdot \vec{M4}}{|\vec{M5}| \cdot |\vec{M4}|} = 0$$

- Kết luận: Top 1 bộ phim khuyến nghị cho user1 là bộ phim có ID là M3.

3. Sử dụng độ đo Jaccard để tìm ra top 2 bộ phim khuyến nghị user1.

- Công thức:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

- Áp dụng:

$$J(M5, M1) = \frac{|M5 \cap M1|}{|M5 \cup M1|} = \frac{0}{11} = 0$$

$$J(M5, M2) = \frac{|M5 \cap M2|}{|M5 \cup M2|} = \frac{1}{14 - 1} = 0.077$$

$$J(M5, M3) = \frac{|M5 \cap M3|}{|M5 \cup M3|} = \frac{3}{12 - 3} = 0.333$$

$$J(M5, M4) = \frac{|M5 \cap M4|}{|M5 \cup M4|} = \frac{0}{12} = 0$$

- Kết luận: Top 2 bộ phim khuyến nghị cho user1 là bộ phim có ID là M3 và M2.