

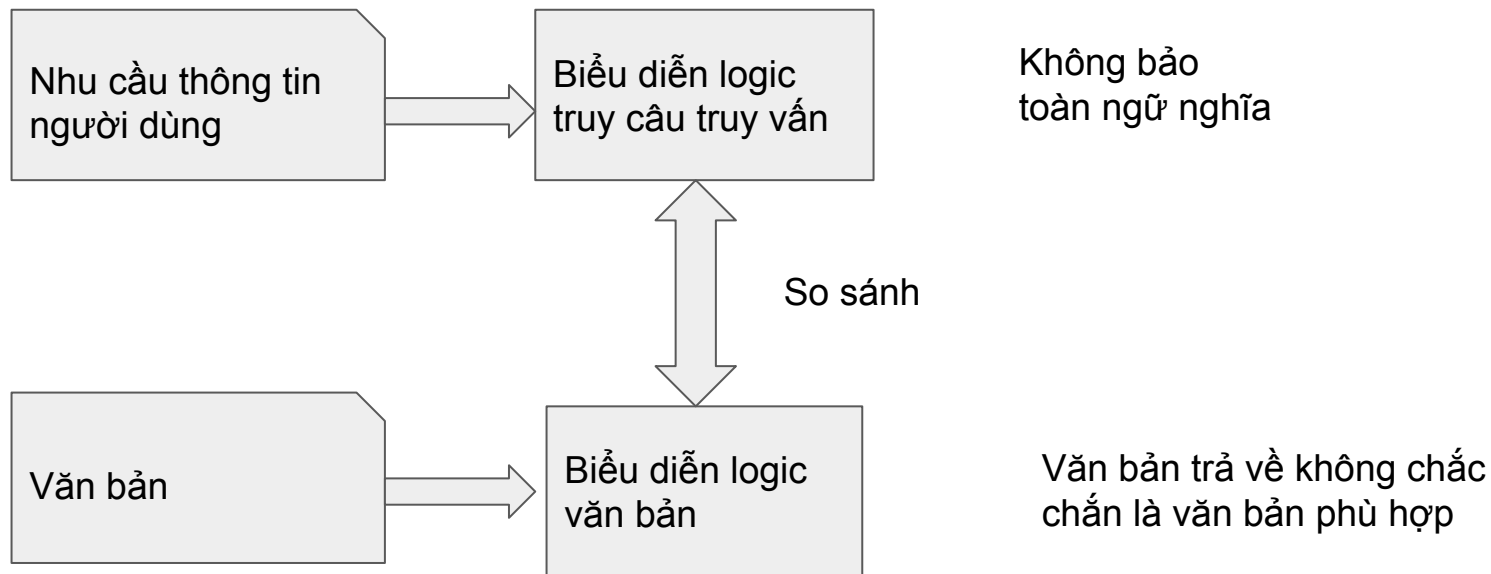
Tìm kiếm và trình diễn thông tin

Mô hình nhị phân độc lập
Probabilistic information retrieval

Nội dung chính

1. Ứng dụng lý thuyết xác suất trong tìm kiếm
2. Mô hình nhị phân độc lập
3. Mô hình (Okapi) BM25

Lý thuyết xác suất trong tìm kiếm thông tin



Có thể ứng dụng lý thuyết xác suất trong tìm kiếm thông tin.

Tổng quan các mô hình xác suất

- Các mô hình xác suất cổ điển:
 - Nguyên tắc xếp hạng xác suất
 - Mô hình nhị phân độc lập,
 - BestMatch25(Okapi)
 - ...
- Tìm kiếm văn bản sử dụng mạng Bayes;
- Các mô hình ngôn ngữ
 - Hướng nghiên cứu mới, hiệu năng cao;

Phương pháp xác suất là một trong những phương pháp đã tồn tại từ lâu nhưng vẫn là đề tài nóng trong tìm kiếm thông tin hiện đại.

Xếp hạng xác suất

Ký hiệu $R_{d,q}$: một biến nhị phân ngẫu nhiên:

$R_{d,q} = 1$ nếu d phù hợp với q ;

$R_{d,q} = 0$, nếu ngược lại.

Theo phương pháp xếp hạng xác suất, các văn bản được trả về theo thứ tự giảm dần giá trị xác suất văn bản phù hợp với truy vấn: $P(R=1|d, q)$.

Trọng số từ

Xếp hạng xác suất: Probabilistic Ranking

“Từ xuất hiện trong những *văn bản đã biết là phù hợp* phải có trọng số cao hơn so với trọng số của từ đó trong *trường hợp không biết những văn bản phù hợp này*.”

“Có thể xây dựng cách tính trọng số từ dựa trên giả thuyết về phân bố từ vựng và luật Bayes.”

Lý thuyết xác suất căn bản (1)

- For events A and B :

$$p(A, B) = p(A \cap B) = p(A | B)p(B) = p(B | A)p(A)$$

- Bayes' Rule

$$p(A | B) = \frac{p(B | A)p(A)}{p(B)} = \frac{p(B | A)p(A)}{\sum_{X=A, \bar{A}} p(B | X)p(X)}$$

Posterior

Prior

- Odds:

$$O(A) = \frac{p(A)}{p(\bar{A})} = \frac{p(A)}{1 - p(A)}$$

Mô hình nhị phân độc lập

Nhị phân: Văn bản được biểu diễn như vector nhị phân đánh dấu sự xuất hiện của từ

- $d = (x_1, \dots, x_n)$
- $x_i = 1$ nếu thuật ngữ thứ i xuất hiện trong d , 0 nếu ngược lại

Độc lập: Sự xuất hiện của mỗi từ trong văn bản là độc lập với những từ còn lại;

Những văn bản khác nhau có thể có cùng một biểu diễn vector.

Mô hình nhị phân độc lập (1)

- Cho truy vấn q
 - Với mỗi văn bản d cần tính xác suất d là tài liệu phù hợp $p(R=1|q, d)$
 - Chỉ quan tâm tới thứ hạng
- Sử dụng cơ hội (Odds) và luật Bayes

- Bayes' Rule

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} = \frac{p(B|A)p(A)}{\sum_{X=A, \bar{A}} p(B|X)p(X)}$$

Posterior

Prior

$$O(A) = \frac{p(A)}{p(\bar{A})} = \frac{p(A)}{1 - p(A)}$$

$$O(R|q, d) = \frac{p(R=1|q, d)}{p(R=0|q, d)} = \frac{\frac{p(R=1|q)p(d|R=1, q)}{p(d|q)}}{\frac{p(R=0|q)p(d|R=0, q)}{p(d|q)}}$$

Mô hình nhị phân độc lập (1)

- Cho truy vấn q
 - Với mỗi văn bản d cần tính xác suất d là tài liệu phù hợp $p(R=1|q, d)$
 - Chỉ quan tâm tới thứ hạng
- Sử dụng cơ hội (Odds) và luật Bayes

$$O(A) = \frac{p(A)}{p(\bar{A})} = \frac{p(A)}{1 - p(A)}$$

$$O(R|q, d) = \frac{p(R=1|q, d)}{p(R=0|q, d)} = \frac{\frac{p(R=1|q)p(d|R=1, q)}{p(d|q)}}{\frac{p(R=0|q)p(d|R=0, q)}{p(d|q)}}$$

$$O(R|q, d) = \frac{p(R=1|q, d)}{p(R=0|q, d)} = \frac{p(R=1|q)}{p(R=0|q)} \cdot \frac{p(d|R=1, q)}{p(d|R=0, q)}$$

Hằng số với
một truy vấn

Cần xác định

Mô hình nhị phân độc lập (2)

$$O(R|q,d) = \frac{p(R=1|q,d)}{p(R=0|q,d)} = \frac{p(R=1|q)}{p(R=0|q)} \cdot \frac{p(d|R=1,q)}{p(d|R=0,q)}$$

Hằng số với
một truy vấn

Cần xác định

Sử dụng giả thuyết độc lập

$$\frac{p(d|R=1,q)}{p(d|R=0,q)} = \prod_{i=1}^n \frac{p(x_i|R=1,q)}{p(x_i|R=0,q)}$$

$$O(R|q,d) = O(R|q) \cdot \prod_{i=1}^n \frac{p(x_i|R=1,q)}{p(x_i|R=0,q)}$$

Mô hình nhị phân độc lập (3)

$$O(R|q,d) = O(R|q) \cdot \prod_{i=1}^n \frac{p(x_i|R=1,q)}{p(x_i|R=0,q)}$$

$$= O(R|q) \cdot \prod_{x_i=1} \frac{p(x_i=1|R=1,q)}{p(x_i=1|R=0,q)} \prod_{x_i=0} \frac{p(x_i=0|R=1,q)}{p(x_i=0|R=0,q)}$$

$$p_i = p(x_i=1|R=1,q)$$

$$1-p_i = p(x_i=0|R=1,q)$$

$$r_i = p(x_i=1|R=0,q)$$

$$1-r_i = p(x_i=0|R=0,q)$$

$$O(R|q,d) = O(R|q) \cdot \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} \cdot \prod_{q_i=1} \frac{1-p_i}{1-r_i}$$

Từ truy vấn có
trong văn bản

Tất cả từ truy₁₈
vấn

Mô hình nhị phân độc lập (5)

$$O(R|q,d) = O(R|q) \cdot \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} \cdot \prod_{q_i=1} \frac{1-p_i}{1-r_i}$$

Hằng số với một truy vấn

Đại lượng duy nhất cần xác định cho mục đích xếp hạng

Hàm xếp hạng

$$\text{Rank}(d,q) = \log \prod_{x_i=q_i=1} \frac{p_i(1-r_i)}{r_i(1-p_i)} = \sum_{x_i=q_i=1} \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

Mô hình nhị phân độc lập (6)

Kết quả tìm kiếm được xác định dựa trên

$$RSV(d, q) = \log \prod_{x_i = q_i = 1} \frac{p_i(1-r_i)}{r_i(1-p_i)} = \sum_{x_i = q_i = 1} \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

$$RSV(d, q) = \sum_{x_i = q_i = 1} c_i; \quad c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

c_i có vai trò như trọng số thuật ngữ trong mô hình này

Những số liệu thống kê cơ bản

Đại lượng thống kê ứng với từ thứ i :

Từ/văn bản	Phù hợp	Không phù hợp	Tổng
$x_i=1$	s	$n-s$	n
$x_i=0$	$S-s$	$N-n-S+s$	$N-n$
Tổng	S	$N-S$	N

$$p_i \approx \frac{s}{S} \quad r_i \approx \frac{n-s}{N-S} \quad p_i = p(x_i=1 | R=1, q); \quad r_i = p(x_i=1 | R=0, q);$$

$$c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)} \quad c_i \approx K(N, n, S, s) = \log \frac{s/(S-s)}{(n-s)/(N-n-S+s)}$$

Làm mịn trọng số

Có thể thêm 0.5 vào mỗi tham số để đảm bảo các trọng số không trở thành vô cùng khi S, s nhỏ:

$$c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)} \quad c_i \approx K(N, n, S, s) = \log \frac{s/(S-s)}{(n-s)/(N-n-S+s)}$$

$$c_t = \log \frac{(s+0.5)(N-S-n+s+0.5)}{(n-s+0.5)(S-s+0.5)}$$

Bắt đầu thực hiện truy vấn

Hoàn toàn không biết về R

$$c_t = \log \frac{N - n + 0.5}{n + 0.5}$$

$$c_t = \log \frac{(s + 0.5)(N - S - n + s + 0.5)}{(n - s + 0.5)(S - s + 0.5)}$$

Tương tự trọng số idf.

Có thể sử dụng giá trị này để tính hạng ban đầu

Bắt đầu thực hiện truy vấn

d Biểu diễn vec-tơ văn bản

	a	b	c	d	e	f	g	h	k	l
1	1			1				1	1	
2								1	1	1
3		1				1	1			
4	1			1						1
5								1	1	
6			1		1					

$$c_t = \log \frac{N - n + 0.5}{n + 0.5}$$

Cải thiện xếp hạng bằng cách ước lượng pi

- Phản hồi từ người dùng
- Sử dụng hằng số $\pi = 0.5$

Cải thiện xếp hạng bằng cách ước lượng pi

Phù hợp phản hồi giả lập

1. Giả sử pi là hằng số với mọi xi trong truy vấn. Ví dụ, $\pi_i = 0.5$ với văn bản bất kỳ
2. Giả sử tập văn bản phù hợp V là tập chứa những văn bản được xếp hạng cao nhất theo mô hình này.
3. Cần xác định lại pi và ri, sử dụng phân bố từ trong V.

Đặt Vi là tập văn bản có chứa xi, chúng ta có

$$\pi_i = (|V_i| + 0.5) / (|V| + 1)$$

4. Giả sử không được trả về đồng nghĩa với không phù hợp,

$$r_i = (n_i - |V_i| + 0.5) / (N - |V| + 1)$$

- 5 Lặp các bước 2-4 cho tới khi hội tụ và trả về kết quả

Ví dụ trọng số phù hợp

d Biểu diễn vec-tơ văn bản

	a	b	c	d	e	f	g	h	k	l
1	1			1				1	1	
2								1	1	1
3		1				1	1			
4	1			1						1
5								1	1	
6			1		1					

R=1

$$c_t = \log \frac{(s+0.5)(N-S-n+s+0.5)}{(n-s+0.5)(S-s+0.5)}$$

Tổng kết mô hình BIM

- Mô hình xác suất dựa trên lý thuyết xác suất để mô hình hóa sự không chắc chắn trong quá trình tìm kiếm
- Sử dụng các giả thuyết về sự độc lập trong quá trình ước lượng giá trị xác suất
- Từ không xuất hiện trong truy vấn không ảnh hưởng tới tính phù hợp (có $p_i = r_i$)
- Trọng số ban đầu của thuật ngữ khi chưa có thông tin về văn bản phù hợp được xác định tương tự idf.
- Phù hợp phản hồi giả lập có thể giúp cải thiện xếp hạng bằng cách xác định lại xác suất thuật ngữ
- Không sử dụng các tần suất thuật ngữ văn bản
- BM25