
Chapitre 6 Modèle probabiliste probabilistic model

1

Modèle probabiliste

- Pourquoi les probabilités ?
 - La RI est un processus incertain et imprécis
 - Imprécision dans l'expression des besoins
 - Incertitude dans la représentation des informations
 - La théorie de la probabilité semble adéquate pour quantifier (pour mesurer) cette incertitude et imprécision

2

Modèle probabiliste

- Le modèle probabiliste tente d'estimer la probabilité qu'un document donné soit pertinent pour une requête donnée
 - $P(\text{pert}/d, q)$: probabilité de pertinence de d vis à vis de q
 - $P(q, d)$
 - $P(q/d)$
 - $P(d/q)$

3

Modèle probabiliste de base

- Considérons une requête q et un document d , le modèle probabiliste tente d'estimer la probabilité que le document d appartienne à la classe des documents pertinents (non pertinents)



$P(\text{NR}|D)$



$P(R|D)$

- Un document est sélectionné si : $P(R/d) > P(\text{NR}/d)$

4

Probability Ranking Principle

- Supposons que la pertinence d'un document est indépendante des autres documents
- Probability Ranking Principle (Principe d'appariement probabiliste)
 - “Ranking documents in decreasing order of probability of relevance to the user who submitted the query, where probabilities are estimated using all available evidence, produces the best possible effectiveness”
- L'efficacité est définie en termes de précision

5

Probability Ranking principal

- Documents triés selon PRP

$$RSV(q,d)=O(d)=P(R/d) / P(NR/d)$$

Rappel Règle de Bayes

$$p(a, b) = p(a \cap b) = p(a | b)p(b) = p(b | a)p(a)$$

$$p(\bar{a} | b)p(b) = p(b | \bar{a})p(\bar{a})$$

$$p(a | b) = \frac{p(b | a)p(a)}{p(b)}$$

6

Probabilistic Ranking Principle Démonstration

- Règle de Bayes

$$p(R | d) = \frac{p(d | R)p(R)}{p(d)}$$

$$p(NR | d) = \frac{p(d | NR)p(NR)}{p(d)}$$

- PRP : Ordonner les documents par rapport

$$RSV(q, d) = O(d) \approx \frac{p(d | R)}{p(d | NR)}$$

7

Probability Ranking principal

- Options :
 - Comment représenter le Document D ?
 - Quelle distribution utilisée pour $P(D | R)$ et $P(D|NR)$?
 - Etant donnée une requête comment estimer les paramètres du modèle ?
- Plusieurs solutions
 - BIR (Binary Independant Model)
 - “Two poisson model”

8

Binary Independence Retrieval (BIR)

- Document : ensemble d'événements
- Événement dénote la présence ou l'absence d'un terme dans un document

$$d = (t_1, \dots, t_n)$$

$$t_i = 1 \quad \text{si un terme est présent dans un document}$$

- **Independence**": les termes apparaissent dans les documents de manière indépendante

9

L'appariement (Probabilistic Ranking Principle)

- Considérons un document comme une liste de termes
 - $P(d/Pert)$ et $P(d/Npert)$ sont estimées par les probabilités conditionnelles selon qu'un terme de la requête est présent dans un document pertinent ou non pertinent.

$$RSV(q, d) = O(d) \approx \frac{p(d(t_1, t_2, \dots, t_n) | R)}{p(d(t_1, t_2, \dots, t_n) | NR)}$$

- en se basant sur l'hypothèse d'indépendance

$$\frac{p(d | R)}{p(d | NR)} = \prod_{i=1}^n \frac{p(t_i | R)}{p(t_i | NR)}$$

10

Binary Independence Model

Retour sur l'hypothèse d'indépendance

- En probabilité, la combinaison de plusieurs événements doit être déterminée comme suit :

$$P(t_1, t_2, t_3, t_4 \dots | R) = P(t_1 | R) * P(t_2 | t_1, R) * P(t_3 | t_1, t_2, R) * P(t_4 | t_1, t_2, t_3, R) * \dots$$

- En RI, la présence et l'absence de termes sont dépendantes.
 - Par exemple, si le terme «informatique» apparaît dans un document, il y a plus de chance que le terme «ordinateur» apparaît aussi. Ainsi

$$P(\text{ordinateur}=1 \mid \text{informatique}=1) > P(\text{ordinateur}=1)$$

11

Binary Independence Retrieval (BIR)

Loi de Bernoulli

$$D = \{t_1=x_1, t_2=x_2, \dots, t_n=x_n\} \quad x_i = \begin{cases} 1 & \text{term present} \\ 0 & \text{term absent} \end{cases}$$

$$P(D | R) = \prod_{i=1}^n P(t_i = x_i | R)$$

$$= \prod_{i=1}^n P(t_i = 1 | R)^{x_i} P(t_i = 0 | R)^{(1-x_i)} = \prod_{i=1}^n p_i^{x_i} (1-p_i)^{(1-x_i)}$$

$$P(D | NR) = \prod_{i=1}^n P(t_i = 1 | NR)^{x_i} P(t_i = 0 | NR)^{(1-x_i)} = \prod_{i=1}^n q_i^{x_i} (1-q_i)^{(1-x_i)}$$

12

Binary Independence Retrieval (BIR)

$$\begin{aligned}
 \text{Odd}(D) &= \log \frac{P(D|R)}{P(D|NR)} = \log \frac{\prod_{i=1}^n p_i^{x_i} (1-p_i)^{(1-x_i)}}{\prod_{i=1}^n q_i^{x_i} (1-q_i)^{(1-x_i)}} \\
 &= \sum_{i:x_i=1}^n x_i \log \frac{p_i(1-q_i)}{q_i(1-p_i)} + \sum_{i=1}^n \log \frac{1-p_i}{1-q_i} \\
 &\propto \sum_{i:x_i=1}^n \log \frac{p_i(1-q_i)}{q_i(1-p_i)}
 \end{aligned}$$

Constante
(quelque
soit le
document)

Comment estimer p_i and q_i

13

Estimation avec des données d'apprentissage

- En considérant pour chaque terme t_i

Documents	Pertinent	Non-Pertinent	Total
$t_i=1$	r	$n-r$	n
$t_i=0$	$R-r$	$N-n-R+r$	$N-n$
Total	R	$N-R$	N

r : nombre de documents pertinents contenant t_i
 n : nombre de documents contenant t_i
 R : nombre total de documents pertinents
 N : nombre de documents dans la collection

14

Estimation par maximum de vraisemblance

Estimation des p_i et q_i

$$p_i = \frac{r}{R} \quad \text{et} \quad q_i = \frac{n-r}{N-R}$$

$$\begin{aligned} RSV(q, d) &= \sum \log \frac{p(1-q)}{q(1-p)} = \\ &= \sum \log \frac{\frac{r}{R} * \frac{N-n-R+r}{N-R}}{\frac{n-r}{N-R} * \frac{R-r}{R}} = \\ &= \sum \log \frac{r/(R-r)}{(n-r)/(N-n-R+r)} \end{aligned}$$

15

Modèle probabiliste BIR

- Lisser les probabilités pour éviter 0,

$$RSV(q, d) = \sum_{t_i \in (d \cap q)} \log \frac{\frac{r_i + 0.5}{R - r_i + 0.5}}{\frac{(n_i - r_i + 0.5)}{(N - n_i - R + r_i + 0.5)}}$$

- Requêtes et documents sont pondérés

$$RSV(q, d_j) = \sum_{t_i \in (d \cap q)} w_{ij} * qtf_i * \log \frac{\frac{r_i + 0.5}{R - r_i + 0.5}}{\frac{(n_i - r_i + 0.5)}{(N - n_i - R + r_i + 0.5)}}$$

- w_{ij} poids du terme i dans le document j
- qtf_i poids du terme i dans la requête q

16

Estimation sans données d'apprentissage

Lorsque des données d'apprentissage pour l'évaluation ne sont pas disponibles

→ estimation a priori: on donne des valeurs pour p_i et q_i

$$p_i = 0.5$$

$q_i = n_i / N$ (l'ensemble des documents non-pertinents est beaucoup plus important que l'ensemble des documents pertinents) revient aussi à considérer qu'on n'a pas d'informations de pertinence dans la

formule

$$RSV(q, d) = \sum_{i \in (q \cap d)} \log\left(\frac{N - n_i}{n_i}\right)$$

Pour éviter le zéro, on ajoute 0,5. RSV devient

$$RSV(q, d) = \sum_{i \in (q \cap d)} \log\left(\frac{N - n_i + 0,5}{n_i + 0,5}\right)$$

17

Pour une Requête et des documents sont pondérés

$$RSV(q, dj) = \sum_{i \in (q \cap d)} W_{ij} * q_{tft} * \log\left(\frac{N - n_i + 0,5}{n_i + 0,5}\right)$$

18

Modèle "2-poisson"

- S. Walker et S. Robertson ont estimé ces paramètres selon la formule : BM25

$$RSV(q,dj) = \sum_{i \in (q \cap dj)} \left(\frac{(k+1) * tfij}{tfij + k * (1 - b + b * (\frac{dlj}{avgdl}))} \right) * \log \left(\frac{(N - ni + 0,5)}{(ni + 0,5)} \right)$$

Tfij : la fréquence du terme *i* dans la document *j*

ni : le nombre de document contenant le terme *i*

N: nombre de document de la collection

dlj : longueur de du document *j* (nombre de termes)

avgdl : longueur moyenne des documents de la collection

k, b : des constantes

19

The BM25 Formula

$$\sum_{T \in Q} w^{(1)} \frac{(k_1 + 1)tf (k_3 + 1)qtj}{K + tf \quad k_3 + qtj}$$

here

Q is a query, containing terms T

$w^{(1)}$ is the Robertson/Sparck Jones weight [5] of T in Q

$$\log \frac{(r + 0.5)/(R - r + 0.5)}{(n - r + 0.5)/(N - n - R + r + 0.5)} \quad (2)$$

N is the number of items (documents) in the collection

n is the number of documents containing the term

R is the number of documents known to be relevant to a specific topic

r is the number of relevant documents containing the term

K is $k_1((1 - b) + b.dl/avgdl)$

k_1, b and k_3 are parameters which depend on the on the nature of the queries and possibly on the database; k_1 and b default to 1.2 and 0.75 respectively, but smaller values of b are sometimes advantageous; in long queries k_3 is often set to 7 or 1000 (effectively infinite)

tf is the frequency of occurrence of the term within a specific document

qtj is the frequency of the term within the topic from which Q was derived

dl and $avgdl$ are respectively the document length and average document length measured in some suitable unit.

Modèle probabiliste : récapitulatif

- Un des modèles les plus importants dans le domaine de la RI
- BM25 donne les meilleures performances en terme de rappel et précision
- Problèmes :
 - Indépendance entre les termes
 - Indépendance entre les documents