

$$RSV_d = \log \sum_{t: x_t=g_t=1} \frac{P_t(1-U_t)}{U_t(1-P_t)} = \sum_{t: x_t=g_t=1} \log \frac{P_t(1-U_t)}{U_t(1-P_t)}$$

	Documents	Relevant	Non-Relevant	TOTAL
Term Present	$x_t=1$	s	$df_t - s$	df_t
Term Absent	$x_t=0$	$S-s$	$(N-df_t) - (S-s)$	$N-df_t$
	TOTAL	S	$N-S$	N

SOL: 含 Q_k 的 document for query $\Rightarrow P_{tk} = P(X_i=1 | R=1, Q_k)$.

$$U_{tk} = P(X_i=1 | R=0, Q_k).$$

假設 P_t, U_t 僅跟 query 中是否有出現該字有關，即：

$$P_t = P_{tk}, U_t = U_{tk}, \forall k: g_{kt}=1.$$

從表格數據可得：

$$P_t = P(X_i=1 | R=1, Q) = \frac{s}{S} \Rightarrow \frac{P_t}{1-P_t} = \frac{s}{S-s}$$

$$U_t = P(X_i=1 | R=0, Q) = \frac{df_t - s}{N-S} \Rightarrow \frac{U_t}{1-U_t} = \frac{df_t - s}{(N-df_t) - (df_t - s)}$$

$$(1) RSV_d = \sum_{t: x_t=g_t=1} \log \frac{P_t(1-U_t)}{U_t(1-P_t)} = \sum_{t: x_t=g_t=1} \log \left(\frac{P_t}{1-P_t} - \log \frac{U_t}{1-U_t} \right)$$

假設 document size 夠大 (即 $N \gg S, N \gg df_t \gg s$)，則

$$\frac{P_t}{1-P_t} = \frac{s}{S} \approx \frac{s}{S}$$

$$\frac{U_t}{1-U_t} = \frac{df_t - s}{(N-df_t) - (df_t - s)} \approx \frac{df_t}{N}$$

$$\Rightarrow RSV_d \approx \sum_{t: x_t=g_t=1} \left(\log \frac{s}{S} - \log \frac{df_t}{N} \right) = \sum_{t: x_t=g_t=1} \left(\log \frac{s}{S} + \boxed{\log \frac{N}{df_t}} \right)$$

IDF.

∴ RSJ 和 TF-IDF 都有考慮 IDF 的資訊，而 RSJ 多考慮了 relevant 的資訊。

TF-IDF 則是多考慮了 TF 的資訊。

$$\log \frac{P_t}{1-P_t}$$

(2) Probabilistic Retrieval Model:

需要有初始假設，亦即假設所有 query term 會在計算出相同的 document relevant value，也就是說，包含相同 terms 的不同文章將產生相同的相關機率值。

Vector Space Model:

不需要任何初始假設，因此在 general cases 下表現較好。

$$\begin{aligned}
 (3) RSV_d &= \sum_{t: x_t=g_t=1} \log \frac{p_t(1-u_t)}{u_t(1-p_t)} \\
 &= \sum_{t: x_t=g_t=1} \log \frac{\frac{s}{N} \cdot \left(1 - \frac{df_t - s}{N-s}\right)}{\frac{d-s}{N-s} \cdot \left(1 - \frac{s}{N-s}\right)} \\
 &= \sum_{t: x_t=g_t=1} \log \frac{\frac{s}{N} \cdot \frac{(N-s)-(df_t-s)}{N-s}}{\frac{df_t-s}{N-s} \cdot \frac{s-s}{N-s}} \\
 &= \sum_{t: x_t=g_t=1} \log \frac{s((N-s)-(df_t-s))}{(df_t-s)(s-s)}
 \end{aligned}$$

$$\begin{aligned}
 \text{以單-term 來討論 ranking} \Rightarrow & \log \frac{s}{s-s} \times \frac{(N-s)-(df_t-s)}{df_t-s} \\
 &= \log \frac{s}{s-s} + \log \frac{(N-s)-(df_t-s)}{df_t-s}
 \end{aligned}$$

分子(s)越大，代表相關文章越容易出現該 query term.

分子(N-s)越大，代表相關文章越不容易出現該 query term.

與 IDF 有相同效果（第 1 小題），故可用來判斷該 query term 是否具有獨特性。

兩個相加的分數越高，代表該 query term 與文章相關。

又具獨特性，因此可用於 ranking。

2. how much wood would a woodchuck chuck
if a woodchuck could chuck wood.
he would chuck he would as much as he could
and chuck as much as a woodchuck would.
if a woodchuck could chuck wood.

sol:

(1) MLE-estimated unigram probability model: $\hat{P}(t|IM_d) = \frac{tf_{t,d}}{|d|}$.

$$\Rightarrow P(\text{would}) = \frac{4}{37}, \quad P(\text{chuck}) = \frac{5}{37}.$$

(2) MLE-estimated bigram model: $\hat{P}(t_1|t_2) = \frac{tf_{t_1,t_2}}{tf_{t_2}}$.

$$\Rightarrow P(\text{wood}|\text{chuck}) = \frac{2}{5}, \quad P(\text{chuck}|\text{would}) = \frac{1}{4}.$$

3. Doc ID Document Text

- 1 click go the shears boy click click click.
- 2 click click.
- 3 metal here
- 4 metal shears click here.

Sol:

$$\text{Mixture model} \Rightarrow P(q|d) \propto \prod_{1 \leq k \leq 8} (\lambda P(t_k|M_d) + (1-\lambda) P(t_k|M_c)), \lambda = \frac{1}{2}$$

	click	go	the	shears	buys	metal	here
Model 1	$\frac{4}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	0	0
Model 2	$\frac{2}{2}$	0	0	0	0	0	0
Model 3	0	0	0	0	0	$\frac{1}{2}$	$\frac{1}{2}$
Model 4	$\frac{1}{4}$	0	0	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{4}$
Model 5	$\frac{7}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{2}{16}$	$\frac{1}{16}$	$\frac{2}{16}$	$\frac{2}{16}$

$$(1) \text{Doc 1} : \frac{1}{2} \times \frac{4}{8} + \frac{1}{2} \times \frac{7}{16} = 0.46875$$

$$\text{Doc 2} : \frac{1}{2} \times \frac{2}{2} + \frac{1}{2} \times \frac{1}{16} = 0.71875$$

$$\text{Doc 3} : \frac{1}{2} \times 0 + \frac{1}{2} \times \frac{1}{16} = 0.21875$$

$$\text{Doc 4} : \frac{1}{2} \times \frac{1}{4} + \frac{1}{2} \times \frac{1}{16} = 0.34375$$

\Rightarrow Ranking: Doc 2 > Doc 1 > Doc 4 > Doc 3.

$$(2) \text{Doc 1} : \frac{1}{2} \times \frac{1}{8} + \frac{1}{2} \times \frac{2}{16} = 0.125$$

$$\text{Doc 2} : \frac{1}{2} \times 0 + \frac{1}{2} \times \frac{2}{16} = 0.0625$$

$$\text{Doc 3} : \frac{1}{2} \times 0 + \frac{1}{2} \times \frac{2}{16} = 0.0625$$

$$\text{Doc 4} : \frac{1}{2} \times \frac{1}{4} + \frac{1}{2} \times \frac{2}{16} = 0.1875$$

\Rightarrow Ranking: Doc 4 > Doc 1 > Doc 2 = Doc 3.

(3). 利用前面 = J. 跟的結果計算 $P(\text{click} | \# \text{doc}) \times P(\text{shears} | \# \text{doc})$

$$\text{Doc 1} : 0.46875 \times 0.125 = 0.0586$$

$$\text{Doc 2} : 0.71875 \times 0.0625 = 0.0449$$

$$\text{Doc 3} : 0.21875 \times 0.0625 = 0.0137$$

$$\text{Doc 4} : 0.34375 \times 0.1875 = 0.0645$$

\Rightarrow Ranking: Doc 4 > Doc 1 > Doc 2 > Doc 3.

4. Doc1: He moved from Taiwan, Taipei to Taiwan, Nantou.

Doc2: He moved from Taiwan, Nantou to Taiwan, Taipei.

Doc3: He moved from Nantou to Taiwan, Taipei.

sol:

$$\text{Mixture model} \Rightarrow P(g_1 | d) \propto \prod_{1 \leq k \leq 18} (\lambda P(t_k | M_d) + (1-\lambda) P(t_k | M_c)), \lambda = \frac{1}{2}$$

$$\text{Query: Taipei Taiwan. } \Rightarrow \text{Doc1: } \left(\frac{1}{2} \times \frac{1}{8} + \frac{1}{2} \times \frac{3}{23}\right) \times \left(\frac{1}{2} \times \frac{2}{8} + \frac{1}{2} \times \frac{5}{23}\right) \approx 0.0298$$

$$\text{Doc2: } \left(\frac{1}{2} \times \frac{1}{8} + \frac{1}{2} \times \frac{3}{23}\right) \times \left(\frac{1}{2} \times \frac{2}{8} + \frac{1}{2} \times \frac{5}{23}\right) \approx 0.0298$$

$$\text{Doc3: } \left(\frac{1}{2} \times \frac{1}{7} + \frac{1}{2} \times \frac{3}{23}\right) \times \left(\frac{1}{2} \times \frac{1}{7} + \frac{1}{2} \times \frac{5}{23}\right) \approx 0.0246$$

$\Rightarrow \text{Ranking: Doc1} = \text{Doc2} > \text{Doc3}$

5.

	DocID	Words in Document	in c = China
Training Set.	1	Taipei Taiwan	Yes
	2	Macao Taiwan Shanghai	Yes
	3	Japan Sapporo	No
	4	Sapporo Osaka Taiwan	No.
Testing Set.	5	Taiwan Taiwan Sapporo	?

sol:

$$(1) P(\text{Taiwan} | c) = \frac{2+1}{5+7} = \frac{3}{12} = \frac{1}{4}$$

$$P(\text{Sapporo} | c) = \frac{0+1}{5+7} = \frac{1}{12}$$

$$P(\text{Taiwan} | \bar{c}) = \frac{1+1}{5+7} = \frac{2}{12} = \frac{1}{6}$$

$$P(\text{Sapporo} | \bar{c}) = \frac{2+1}{5+7} = \frac{3}{12} = \frac{1}{4}$$

$$(2) P(c | D_5) = \frac{3}{4} \left(\left(\frac{1}{4}\right)^2 \cdot \frac{1}{12} \right) \approx 0.0039 \Rightarrow \text{Yes!!}$$

$$P(\bar{c} | D_5) = \frac{1}{4} \left(\left(\frac{1}{6}\right)^2 \cdot \frac{1}{4} \right) \approx 0.0019$$

$$(3) P(\text{Taipei} | c) = P(\text{Macao} | c) = P(\text{Shanghai} | c) = \frac{1+1}{2+2} = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{Taiwan} | c) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$P(\text{Japan} | c) = P(\text{Osaka} | c) = P(\text{Sapporo} | c) = \frac{0+1}{2+2} = \frac{1}{4}$$

$$P(\text{Taipei} | \bar{c}) = P(\text{Macao} | \bar{c}) = P(\text{Shanghai} | \bar{c}) = \frac{0+1}{2+2} = \frac{1}{4}$$

$$P(\text{Taiwan} | \bar{c}) = P(\text{Japan} | \bar{c}) = P(\text{Osaka} | \bar{c}) = \frac{1+1}{2+2} = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{Sapporo} | \bar{c}) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$\begin{aligned}
 (4) \cdot P(c|D_S) &\propto P(c) \cdot P(\text{Taiwan}|c) \cdot P(\text{Sapporo}|c) \cdot (1 - P(\text{Taipei}|c)) \cdot (1 - P(\text{Macau}|c)) \\
 &\quad \cdot (1 - P(\text{Shanghai}|c)) \cdot (1 - P(\text{Japan}|c)) \cdot (1 - P(\text{Osaka}|c)) \\
 &= \frac{3}{4} \times \frac{3}{4} \times \frac{1}{4} \times (1 - \frac{1}{2}) \times (1 - \frac{1}{2}) \times (1 - \frac{1}{2}) \times (1 - \frac{1}{4}) \times (1 - \frac{1}{2}) \\
 &\doteq 0.0099
 \end{aligned}$$

$$\begin{aligned}
 P(c|D_S) &\propto P(c) \cdot P(\text{Taiwan}|c) \cdot P(\text{Sapporo}|c) \cdot (1 - P(\text{Taipei}|c)) \cdot (1 - P(\text{Macau}|c)) \\
 &\quad \cdot (1 - P(\text{Shanghai}|c)) \cdot (1 - P(\text{Japan}|c)) \cdot (1 - P(\text{Osaka}|c)) \\
 &= \frac{1}{4} \times \frac{1}{2} \times \frac{3}{4} \times (1 - \frac{1}{2}) \times (1 - \frac{1}{4}) \times (1 - \frac{1}{2}) \times (1 - \frac{1}{2}) \\
 &\doteq 0.0099
 \end{aligned}$$

6. (經典問題太長跳過了~)

SOL:

由 RST 模型，使用 $P(R=1|D, Q)$ 做 ranking 等價於使用 odds $O(R|D, Q)$ 做 ranking, 又。

$$O(R|D, Q) = O(R|D) \cdot \frac{P(D|R=1, Q)}{P(D|R=0, Q)}$$

其中 $O(R|D)$ 與 document 無關可忽略，故 ranking 等價於 $\frac{P(D|R=1, Q)}{P(D|R=0, Q)}$

$$\text{使用 multinomial 的假設} \Rightarrow P(D|R=1, Q) = \frac{n!}{x_1! x_2! \dots x_v!} p_1^{x_1} p_2^{x_2} \dots p_v^{x_v}$$

$$P(D|R=0, Q) = \frac{n!}{x_1! x_2! \dots x_v!} u_1^{x_1} u_2^{x_2} \dots u_v^{x_v}$$

其中 $x_i = c(w_i, D)$, $p_i = P(W=w_i|R=1, Q)$, $u_i = P(W=w_i|R=0, Q)$.

故 ranking 亦可以寫成：

$$\frac{P(D|R=1, Q)}{P(D|R=0, Q)} = \frac{\frac{n!}{x_1! x_2! \dots x_v!} p_1^{x_1} p_2^{x_2} \dots p_v^{x_v}}{\frac{n!}{x_1! x_2! \dots x_v!} u_1^{x_1} u_2^{x_2} \dots u_v^{x_v}} = \prod_{i=1}^v \frac{p_i^{x_i}}{u_i^{x_i}}$$

$$\text{取 log} \Rightarrow \log \prod_{i=1}^v \frac{p_i^{x_i}}{u_i^{x_i}} = \sum_{i=1}^v \log \left(\frac{p_i}{u_i} \right)^{x_i}$$

$$= \sum_{i=1}^v x_i \log \frac{p_i}{u_i}$$

$$= \sum_{w \in V} c(w_i, D) \log \frac{P(W=w_i|R=1, Q)}{P(W=w_i|R=0, Q)}$$

因此，需要估計的參數為所有 p_i 和 u_i ，共 $= |V|$ 個。

但 $\sum_{i=1}^v p_i = \sum_{i=1}^v u_i = 1$ 為定值，故 free parameters 只 $(|V|-2)$ 個。

(2) 已知 D_j 與 query 及無關，因此估計出來的 $P(D_j|R=0, \theta)$ 應盡可能大。

套用 MLE 去估計 $u_i = P(W_i|R=0, \theta)$:

$$\operatorname{argmax}_{u_i} \prod_j P(D_j|R=0, \theta) = \operatorname{argmax}_{u_i} \sum_j \log P(D_j|R=0, \theta).$$

$$\text{令 } x_{ij} = c(w_i, D_j), x_{1j} + x_{2j} + \dots + x_{nj} = m_j = |D_j|.$$

$$\text{使用 multinomial 的假設} \Rightarrow P(D_j|R=0, \theta) = \frac{m_j!}{\prod_i x_{ij}!} \prod_i u_i^{x_{ij}}.$$

可將問題變成：

$$\begin{aligned} \operatorname{argmax}_{u_i} \sum_j \log P(D_j|R=0, \theta) &= \operatorname{argmax}_{u_i} \sum_j \log \frac{m_j!}{\prod_i x_{ij}!} \prod_i u_i^{x_{ij}} \\ &= \operatorname{argmax}_{u_i} \sum_j (\log m_j! - \sum_i \log x_{ij}! + \sum_i x_{ij} \log u_i) \end{aligned}$$

由於每個字都是整個 vocabulary 中的其中一個字，故 $\sum_i u_i = \sum_i P(W_i|R=0, \theta) = 1$.

機率總和為 1.

套用 Lagrange multiplier，為了尋找多元函數在某變數受到一個或多個條件的約束時的極值，將一個有 n 個變數與 k 個約束條件的最佳化問題轉換成一個有 $(n+k)$ 個變數的方程組的解的問題：

$$L = \sum_j (\log m_j! - \sum_i \log x_{ij}! + \sum_i x_{ij} \log u_i) + \lambda (1 - \sum_i u_i)$$

各項的偏微分皆為 0：

$$\frac{\partial L}{\partial u_i} = \sum_j \frac{x_{ij}}{u_i} - \lambda = 0 \Rightarrow \sum_j x_{ij} = \lambda u_i \quad \text{--- ①}$$

$$\frac{\partial L}{\partial \lambda} = 1 - \sum_i u_i = 0 \Rightarrow \sum_i u_i = 1 \Rightarrow \lambda = \sum_i u_i \quad \text{--- ②}$$

將 ①、② 合併整理化簡後可得：

$$\lambda = \sum_i \lambda u_i = \sum_i \sum_j x_{ij} = \sum_j \sum_i x_{ij} = \sum_j |D_j|.$$

$$\Rightarrow u_i = \frac{\sum_j x_{ij}}{\lambda} = \frac{\sum_j c(w_i, D_j)}{\sum_j |D_j|}.$$

$$\text{故, } P(W|R=0, \theta) = \frac{\sum_j c(w_i, D_j)}{\sum_j |D_j|}.$$

(3) 令 $D_Q = Q \rightarrow$ 視為 the only example of a relevant document.

因此，估計出來的 $P(D_Q | R=1, Q)$ 應盡可能大。

套用 MLE 去估計 $P_i = P(W_i | R=1, Q)$:

$$\operatorname{argmax}_{P_i} P(D_Q | R=1, Q) = \operatorname{argmax}_{P_i} \log P(D_Q | R=1, Q).$$

令 $X_i = c(W_i, D_Q)$, $X_1 + X_2 + \dots + X_v = m = |D_Q|$.

$$\text{使用 multinomial 的假設} \Rightarrow P(D_Q | R=1, Q) = \frac{m!}{\prod_i X_i!} \prod_i P_i^{X_i}.$$

可將問題寫成：

$$\begin{aligned} \operatorname{argmax}_{P_i} \log P(D_Q | R=1, Q) &= \operatorname{argmax}_{P_i} \log \frac{m!}{\prod_i X_i!} \prod_i P_i^{X_i} \\ &= \operatorname{argmax}_{P_i} \log m! - \sum_i \log X_i! + \sum_i X_i \log P_i \end{aligned}$$

由於每個字都是整個 vocabulary 中的其中一個字，故 $\sum_i P_i = \sum_i P(W_i | R=1, Q) = 1$.

概率總和為 1.

套用 Lagrange multiplier，為了尋找多元函數在具變數受到一個或多個條件的約束時的極值，將一個有 n 個變數和 k 個約束條件的最佳化問題轉換成一個有 $(n+k)$ 個變數的方程組的解的問題：

$$L = (\log m! - \sum_i \log X_i! + \sum_i X_i \log P_i) + \lambda (1 - \sum_i P_i)$$

各項的偏微分必須為 0：

$$\frac{\partial L}{\partial P_i} = \frac{X_i}{P_i} - \lambda = 0 \Rightarrow X_i = \lambda P_i \quad \text{③}$$

$$\frac{\partial L}{\partial \lambda} = 1 - \sum_i P_i \Rightarrow \sum_i P_i = 1 \Rightarrow \lambda = \sum_i \lambda P_i \quad \text{④}$$

將 ③、④ 合併整理化簡後可得：

$$\lambda = \sum_i \lambda P_i = \sum_i X_i = |D_Q|.$$

$$\Rightarrow P_i = \frac{X_i}{\lambda} = \frac{c(W_i, D_Q)}{|D_Q|}$$

$$\text{故, } P(W_i | R=1, Q) = \frac{c(W_i, Q)}{|Q|}.$$

(4) 套用以上結果：

$$\begin{aligned} \text{Score}(Q, D) &= \sum_{w \in Q} c(w, D) \log \frac{P(w|R=1, Q)}{P(w|R=0, Q)} \\ &= \sum_{w \in Q} c(w, D) \left[\log \frac{c(w, Q)}{|Q|} - \log \frac{\sum_{j=1}^n c(w, D_j)}{\sum_{j=1}^n |D_j|} \right] \end{aligned}$$

① 由於 $c(w, Q)$ 在算詞頻，故該項就像 TF 的方法

② 由於 $\sum_{j=1}^n |D_j|$ 在算 collection 的大小，而 $\sum_{j=1}^n c(w, D_j)$ 在算 word 在 collection 中的數量，兩項的比值就類似 IDF 的算法（皆為 word 數量除以 collection 大小的比值取 \log ）。

故，此方法有 TF 跟 IDF 的效果，但由於 score 中沒有根據 document 的大小做修正，故沒有 document length normalization 的功能。

(5) 因為題目的式子未考慮文章的長度，故效果會比有考慮該問題的 BM25 差。

方括號中第一項 $c(w, Q)$ 為 0 時，會出現 $\log 0 = -\infty$ ，故可在分子加上 $\frac{1}{2}$ 。另外由於分母是由 $|Q| = \sum_{w \in Q} c(w, Q)$ 計算來的，故須同時在分母加上 $\frac{|V|}{2}$ ；同樣，也可以在第二項分子加上 $\frac{1}{2}$ ，並在分母加上 $\frac{|V|}{2}$ ，以解決 zero probability problem。

另外，若將整個式子除以 document 的長度 $|D|$ ，則會使得較長的 document 有較多的調整，進而達到 document length normalization 的效果。

套用以上 trick 後，新的 Score 為：

$$\hat{\text{Score}}(Q, D) = \sum_{w \in Q} \frac{c(w, D)}{|D|} \left[\log \frac{c(w, D) + \frac{1}{2}}{|Q| + \frac{|V|}{2}} - \log \frac{\sum_{j=1}^n c(w, D_j) + \frac{1}{2}}{\sum_{j=1}^n |D_j| + \frac{|V|}{2}} \right]$$

證：在這樣的 normalization 效果下，同一篇 document 重複兩次所得到的新 document 的分數會一樣，也就是說，假設 $\bar{D} = D + D$ ，由於 $c(w, \bar{D}) = 2c(w, D)$ ， $|\bar{D}| = 2|D|$ 有相同的比例關係，而方括號內的項與對於 \bar{D} 和 D 是相同的，因此可以推得 $\hat{\text{Score}}(Q, \bar{D}) = \hat{\text{Score}}(Q, D)$ ，可避免長度較長的文章得到不應得的高分數。