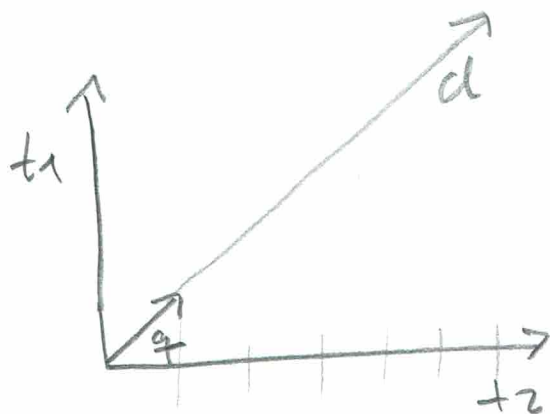


RI- CHAPTER Retrieval Models



$$g = t_1 t_2$$

$$d = t_1 t_1 t_1 t_1 t_1 t_1 t_2 t_2 t_2 t_2 t_2 t_2$$

\vec{g} y \vec{d} hablan de lo mismo pero de manera verbosa. Distancia Euclidea de \vec{g} y \vec{d} es grande, intuitivamente no es una buena medida para lo que pretendemos

$$\vec{d} \cdot \vec{g} = |\vec{d}| \cdot |\vec{g}| \cdot \cos(\angle d, g) \quad \text{Español técnico con términos}$$

$$\cos(\angle d, g) = \frac{\vec{d} \cdot \vec{g}}{|\vec{d}| \cdot |\vec{g}|} = \frac{\sum_{j=1}^t d_j \cdot g_j}{\sqrt{\sum_{j=1}^t d_j^2} \sqrt{\sum_{j=1}^t g_j^2}} \quad \text{rank}$$

$$\text{rank} = \frac{\sum_{j=1}^t d_j \cdot g_j}{\sqrt{\sum_{j=1}^t d_j^2}} = \frac{\vec{d} \cdot \vec{g}}{|\vec{d}|}$$

- term frequency para el término k en el documento i

$$tf_{ik}$$

- nº ocurrencias del término k en el documento i

$$f_{ik} \quad \text{raw } tf$$

- Heurística

$$tf \text{ normalizado} \quad tf_{ik} = \frac{f_{ik}}{\sum_{j=1}^t f_{ij}}$$

Suma de las ocurrencias de todos los términos del documento i

$$tf \text{ logarítmico} \quad tf_{ik} = \begin{cases} 0 & \text{si } f_{ik} = 0 \\ 1 + \log f_{ik} & \text{si } f_{ik} \geq 1 \end{cases}$$

- Idf, inverse document frequency para el término k . Idf logarítmico:

$$idf_k = \log \frac{N}{n_k}$$

N , num. docs de la colección

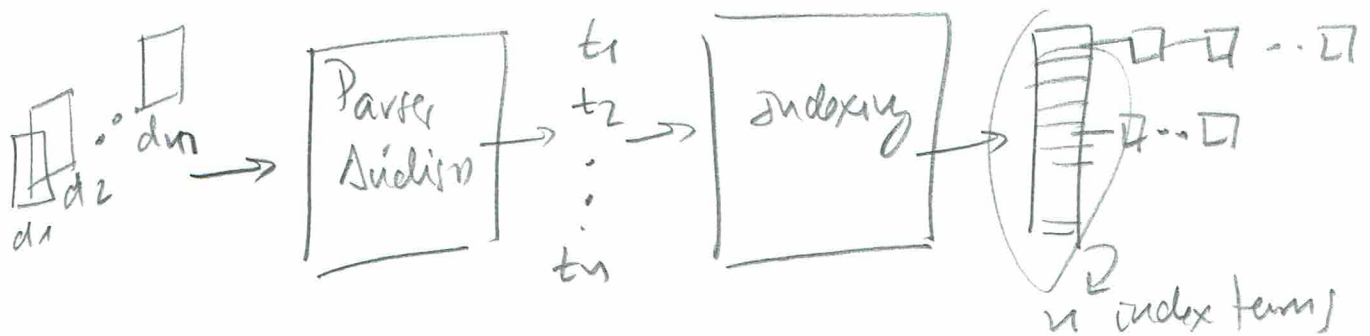
n_k , num docs en las que aparece el término k

$$\text{Si } n_k = N \rightarrow idf_k = 0$$

$$\text{Si } n_k = 1 \rightarrow idf_k = \log N$$

VSIM Vector Space Model

m docs
n terms



Docs y Queries son vectores de dimension n

$$d = \langle d_1, d_2, \dots, d_n \rangle$$

$$q = \langle q_1, q_2, \dots, q_n \rangle$$

$$\cos(q, d) = \frac{\sum_{i=1}^n q_i d_i}{|q| |d|} = \frac{\sum_{i=1}^n q_i d_i}{\sqrt{\sum_{i=1}^n q_i^2} \sqrt{\sum_{i=1}^n d_i^2}} = \text{rank}$$

$$\text{rank} = \frac{\sum_{i=1}^n q_i d_i}{\sqrt{\sum_{i=1}^n d_i^2}} = \sum_i q_i \frac{d_i}{\sqrt{\sum_{i=1}^n d_i^2}}$$

componente
vector
normalizado

$$\frac{\sum_{i=1}^n q_i d_i}{\sqrt{\sum_{i=1}^n d_i^2}} \rightarrow \text{components de } \vec{q} \text{ y } \vec{d} \text{ normalizados}$$

The diagram shows two overlapping circles representing the vectors \vec{q} and \vec{d} . The left circle contains the expression $\frac{\sum_{i=1}^n q_i}{\sqrt{\sum_{i=1}^n q_i^2}}$ and the right circle contains the expression $\frac{\sum_{i=1}^n d_i}{\sqrt{\sum_{i=1}^n d_i^2}}$. An arrow points from the circles to the text 'components de \vec{q} y \vec{d} normalizados'.

Equation pseudo

$g_i, d_i \rightarrow$ binario 1/0

$\rightarrow tf$ raw tf del término en $\rightarrow d$
 $\rightarrow g$

$\rightarrow tf$ normalizado
en el doc
o en la query

$$\frac{tf_i}{\sum_{k=1}^n tf_k}$$

$\rightarrow tf \log (1 + \log tf)$

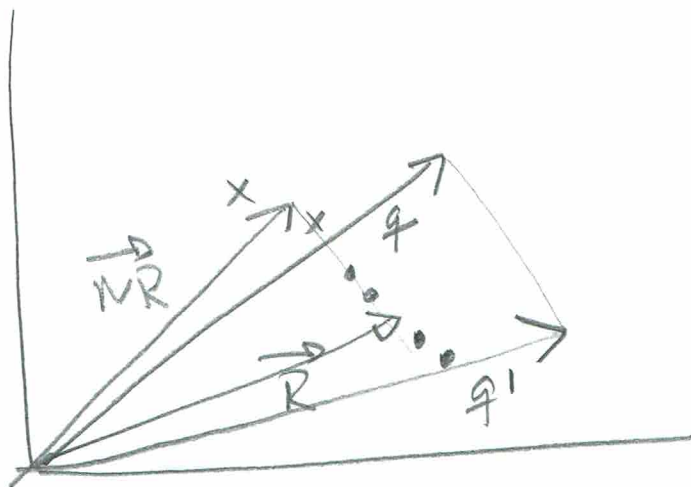
$\rightarrow \underline{edf} = \log \frac{N}{n_i}$

$\boxed{tf \cdot edf}$

\rightarrow perc code
 g_i, d_i

Relevance Feedback

- docs relevant
- x doc, no relevant



$$\vec{q}' = \vec{q} + \vec{R} - \vec{NR}$$

\vec{R} = centroide de los docs relevantes

\vec{NR} = centroide de los docs no relevantes

Con parámetros α, β, γ

$$\vec{q}' = \alpha \vec{q} + \beta \vec{R} - \gamma \vec{NR}$$

En la práctica puede llevar a:

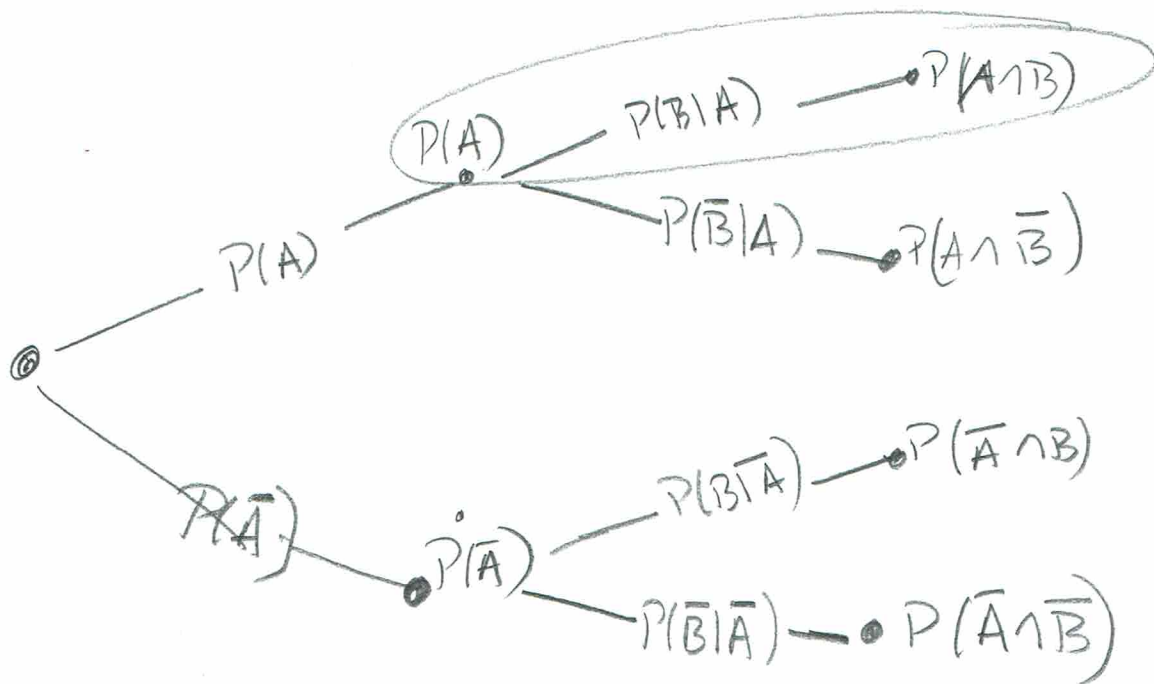
1. Componentes negativas en \vec{q}'
2. \vec{q}' con muchos términos no nulos

soluciones

1. Ignorarlos

2. Seleccionar términos con mayor tf-idf

Prob. condicional



$$P(A) \cdot P(B|A) = P(A \cap B)$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Queremos estimar $P(d|q)$, para todos los docs d
 dado una query q , es decir, queremos
 estimar

$P(d_1|q)$, $P(d_2|q)$, $P(d_3|q)$, etc

y producir ranking por orden decreciente
 de $P(d|q)$

$$P(d|q) = \frac{P(d \cap q)}{P(q)}$$

$$P(q|d) = \frac{P(q \cap d)}{P(d)}$$

$$\rightarrow \boxed{P(d|q) = \frac{P(q|d) \cdot P(d)}{P(q)}} \quad \text{rank}$$

$$\text{rank} = \boxed{P(q|d) \cdot P(d)}$$

→ QUERY LIKELIHOOD

↪ q tiene q_i terms, suponiendo los $q_i|D$
 son eventos independientes, modelo
 multinomial

$$\boxed{P(q|d) = \prod_{i=1}^n P(q_i|D)}$$

q time n query terms

Si a y b son independientes

$$P(a \cap b) = P(a|b) = P(a) \cdot P(b)$$

$$P(a|b) = \frac{P(a \cap b)}{P(b)} = \frac{P(a) \cdot P(b)}{P(b)}$$

→ si son independientes

$$P(a|b) = P(a), \quad y \quad P(b|a) = P(b)$$

Lo mismo para probs condicionales. Si a y b son condicionalmente independientes de c :

$$P(a, b | c) = P(a | c) \cdot P(b | c)$$

Esto fue lo que supusimos al estimar

el QL

$$P(q | d) = P(q_1, q_2, \dots, q_n | d) = \prod_{i=1}^n P(q_i | d)$$

q tiene n términos. Estos pueden ser t términos distintos en q , que cada uno aparece t_f veces

$$P(q | d) = \prod_{i=1}^n P(q_i | d) = \prod_{\substack{\text{términos} \\ \text{distintos} \\ t \text{ en } q}} P(t | d)^{t_f \cdot t, q}$$

average likelihood: MODELO MULTINOMIAL
(omitendo factor de la distribución multinomial).

Distribución binomial con prob. p , parámetros n, K

Una variable aleatoria X se rige por esta distribución si

$$P(K; n, p) = P(X=K) = \binom{n}{K} p^K (1-p)^{n-K}$$

Representa la probabilidad de que ocurran K éxitos en n intentos (n sucesos de Bernoulli, i.e., sucesos binarios, independientes con prob. p)

Por lo tanto: la probabilidad de éxito en un suceso Bernoulli, p

K éxitos ocurren con prob. p^K
 $n-K$ fallos " " " " $(1-p)^{n-K}$

Además los K éxitos pueden ocurrir en cualquier lugar de los n intentos, lo que implica que hay $\binom{n}{K} = \frac{n!}{K!(n-K)!}$

formas distintas de distribuir los K éxitos.

Ej: Prob. de sacar 3 unos en 10 intentos con un dado de 6 caras

$$P(3; 10, \frac{1}{6}) = \binom{10}{3} \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^7 = \underline{0.155}$$

La distribución multinomial es una generalización de la binomial, donde cada intento lleva al éxito por K distintas categorías

Distribución multinomial

Se extraen n bolas de una urna con bolas de K diferentes colores, reemplazando la bola extraída después de cada extracción.

$X_i \rightarrow$ var. para el número de bolas extraídas del color i ($i=1, \dots, K$)

$p_i \rightarrow$ prob. de la bola extraída sea del color i

Prob. de extraer x_i bolas de cada color i para los K diferentes colores en n extracciones con reemplazo:

$$f(x_1, \dots, x_K; n, p_1, \dots, p_K) = P(X_1=x_1, X_2=x_2, \dots, X_K=x_K)$$

$$= \begin{cases} \frac{n!}{x_1! x_2! \dots x_K!} p_1^{x_1} \cdot p_2^{x_2} \cdot \dots \cdot p_K^{x_K} & \text{si } \sum_{i=1}^K x_i = n \\ 0 & \text{otherwise} \end{cases}$$

factor de la distribución multinomial

Analogía con QL

urna \rightarrow documento

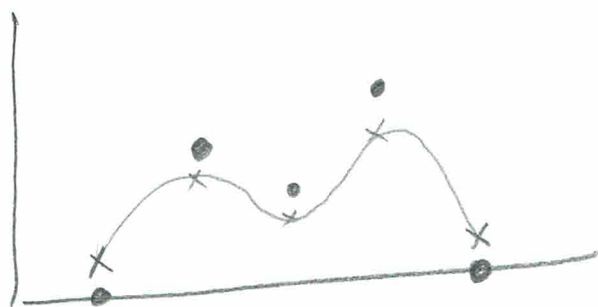
n , num. bolas extraídas \rightarrow n query terms totales

K diferentes colores \rightarrow K diferentes palabras (n palabras únicas)

$$p_i = \frac{n^o \text{ bolas color } i}{n^o \text{ bolas totales urna}} \rightarrow p(w|d) = \frac{n^o \text{ ocurrencias } w \text{ en } d}{n^o \text{ ocurrencias total de palabras en } d}$$

Smoothing

- Probs discretas originales
- x Probs suavizadas



$$P(q_i|C) = \frac{\text{n\u00b0 de ocurrencias palabra } q_i \text{ en } C}{\text{n\u00b0 de ocurrencias totales en } C}$$

$$P(Q|D) \leadsto \log P(Q|D)$$

\log es transformaci\u00f3n mon\u00f3tona

$$\text{si } A > B \Rightarrow \log(A) > \log(B) \Rightarrow$$

\Rightarrow preserva el ranking

Slide where is tf-idf weight?

$i = 1, \dots, n$ todos los query terms

$i: f(q_i, D) > 0$

$i: f(q_i, D) = 0$

terminos de la query que ocurre en el doc.
" " " " " " " " " " " "

De la 2\u00aa a 3\u00aa linea, valemos $\frac{1}{\sum_{i: f(q_i, D) > 0} \log \left(\frac{1/c q_i}{|C|} \right)}$

De la 3\u00aa a 4\u00aa linea, $\left(\sum_{i=1}^n \log \left(\frac{1/c q_i}{|C|} \right) \right)$ no afecta al ranking, igual por todos los docs.

Se le da a algo parecido a tf-idf pero con una derivaci\u00f3n formal

Dirichlet Smoothing

$$\alpha_D = \frac{\mu}{|D| + \mu}$$

$$P(\underline{q_i | D}) = (1 - \alpha_D) \frac{f_{q_i | D}}{|D|} + \alpha_D \frac{c_{q_i}}{|C|} =$$

$$= \frac{\cancel{|D|}}{|D| + \mu} \cdot \frac{f_{q_i | D}}{\cancel{|D|}} + \frac{\mu}{|D| + \mu} \cdot \frac{c_{q_i}}{|C|} =$$

$$= \frac{|C| f_{q_i | D} + \mu \frac{c_{q_i}}{|C|}}{(|D| + \mu) \cdot |C|} = \frac{f_{q_i | D} + \mu \frac{c_{q_i}}{|C|}}{|D| + \mu}$$

$$P(Q | D) = \prod_{i=1}^n P(q_i | D) \iff QL$$

$$\log P(Q | D) = \sum_{i=1}^n \log \frac{f_{q_i | D} + \mu \frac{c_{q_i}}{|C|}}{|D| + \mu}$$

QUERY LIKELIHOOD MODEL

$$P(Q|D) = \prod_{q_i} P(q_i|D)$$

q_i query term

DOCUMENT LIKELIHOOD MODEL

$$P(D|Q) = \prod_{d_i} P(d_i|Q) \text{ } d_i \text{ document words}$$

→ • dificultad de comparar Docs de muy distinta longitud

$$\bullet | \{d_i\} | \Rightarrow | \{q_i\} |$$

aunque muchos d_i no están en Q , si no están hay que calcular igual $P(d_i|Q)$

RELEVANCE MODEL

$$P(w|R)$$

KL-DIVERGENCE (Kullback-Leibler)

Mide distancias entre distribuciones de prob.

Si $Q(x)=0 \geq P(x) \neq 0$, no está definido, por lo que si $Q(x)=0$, la conclusión es $KL-DIV \geq 0$

$$KL(P||Q) = \sum_x P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

Siempre no negativa

$P(w|R) \rightsquigarrow P$

$P(w|D) \rightsquigarrow Q$, una aproximación de P

valores pequeños de $KL \rightarrow$ menor distancia \rightarrow mayor similitud $\rightarrow -KL$ nos da directamente el

ranking por similitud

$$\begin{aligned} -KL(P(w|R) || P(w|D)) &= -\sum_{w \in V} P(w|R) \log \frac{P(w|R)}{P(w|D)} = \\ &= \sum_{w \in V} P(w|R) \log P(w|D) - \sum_{w \in V} P(w|R) \log P(w|R) \end{aligned}$$

igual para todos los docs \rightarrow no afecta ranking

Selección KL-Div y QL

$$QL \quad P(Q|D) = \prod_{i=1}^n P(q_i|D) \stackrel{\text{rank } n}{=} \sum_{i=1}^n \log P(q_i|D)$$

Si tenemos $P(w|R)$ con $\frac{f_{w,Q}}{|Q|}$ (MLE)

$$KL Div = \sum_{w \in V} P(w|R) \log P(w|D) =$$

$$= \sum_{w \in V} \left(\frac{f_{w,Q}}{|Q|} \right) \log P(w|D)$$

0 Para las w que no están en Q \rightarrow

\rightarrow producen el mismo ranking

Estimación de RM

Definimos $P(w|R) \approx P(w|q_1, q_2, \dots, q_n)$

Por definición de probabilidad condicional:

$$P(w|R) \approx \frac{P(w, q_1, q_2, \dots, q_n)}{P(q_1, q_2, \dots, q_n)}$$

Eliminación de RM

h DEC y partición espacio eventos

ley prob. total

• ley Prob. Total.

$$P(w, q_1, \dots, q_n) = \sum_{DEC} P(D) P(w, q_1, \dots, q_n | D)$$

• Ahora suponemos que w y las query wrds q_i son condicionalmente independientes de D , i.e., una vez elegido D , w y las q_i 'are i.i.d. (identically and independently) sampled'



$$P(w, q_1, \dots, q_n | D) = P(w | D) \cdot \prod_{i=1}^n P(q_i | D)$$

$$P(w, q_1, \dots, q_n) = \sum_{DEC} P(D) P(w | D) \prod_{i=1}^n P(q_i | D)$$

$$P(w | R) = \frac{P(w, q_1, \dots, q_n)}{P(q_1, \dots, q_n)} = \frac{P(w, q_1, \dots, q_n)}{\sum_{w \in V} P(w, q_1, \dots, q_n)}$$

↙ ↘
cte de normalización

Pseudo Feedback Alg

En 4. A efectos prácticos seleccionamos las palabras con $P(w|Z)$ mas alto, y el ranking que produce 4 es el mismo que lanzando una query con QZ con estas top w queries terms

RM3 $P'(w|Z) = \lambda P(w|Z) + (1-\lambda) P(w|Q)$
↳ Interpola el RM con la query original.

Eval de PRF Alg, además de
MAP, NDCG, P@k, Robustness Index

$$R_I = \frac{n_+ - n_-}{|Q|}$$

n_+ n_- queries mejoradas por PRF Alg
 n_- " " dadas por " "

Line PRF with Linear Regression Methods
ACM SIGIR

Line > Relevance Models
For reference