

23-08-24

TF-IDF

#score of a doc for a given query.

$$\text{score}(q, d) = \sum_{t \in q \cap D} \text{tf} \cdot \text{idf}_t$$

where

tf = term frequency

idf = inverse doc. freq.

Document as vector

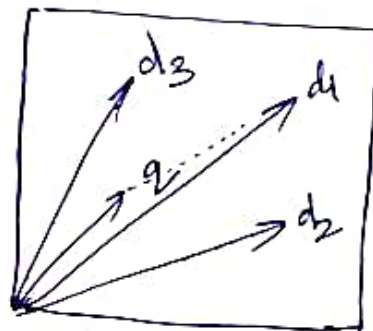
terms - axis/axes

doc - points in space

query - vectors in space

similarity - Rank the doc. based on the proximity.

mag + direction



$$\text{distance} \propto \frac{1}{\text{similarity}}$$

28-08-24

EUCLIDIAN DISTANCE

Distance b/w (x_1, y_1) & (x_2, y_2)

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

- Drawback : distance will not be the only parameter to find out the similarity.
∵ large distance can also have similarity.

Cosine Similarity b/w 3 docs

term	SaS	PaP	WH
affection	115	58	20
jealous	10	7	11
gossip	2	0	6
wuthering	0	0	38

Abbr

SaS : Sense & Sensibility } Jane Austen
PaP : Pride & Prejudice }
WH : Wuthering Height - Emily Bronte

23-08-24

$$W = 1 - \log_{10}$$

i) Log. freq. weighing

term	SaS	PaF	W11
Affection	3.06	2.76	2.3
jealous	2	1.84	2.04
Gossip	1.3	0	1.77
withering	0	0	2.57

ii) cosine similarity of length normalized vector

$$\cos(\vec{q}, \vec{x}) = \frac{\vec{q} \cdot \vec{x}}{|\vec{q}| |\vec{x}|}$$

$$= \frac{1}{\sqrt{N}} \sum_{i=1}^N q_i x_i$$

where,

N = no. of terms

$N = 4$. (in above example).

iii) To normalize

$$\frac{w_i}{\sqrt{\sum w_i^2}}$$

$$\text{eg: } \frac{3.06}{\sqrt{3.06^2 + 2^2 + 1.3^2}} = \frac{3.06}{\sqrt{15.0636}} = \frac{3.06}{3.88}$$

$$\sqrt{19.1894} = 4.38$$

$$\sqrt{11.0032} = 3.32$$

ii) length-normalized vector

term	$w_i/3.87$ SaS	$w_i/3.32$ PaP	$w_i/4.38$ WH
Affection	0.789	0.83	0.525 $= 0.53$
jealous	0.517 $= 0.52$	0.55	0.466 $= 0.47$
gossip	0.3136 $= 0.34$	0	0.404 $= 0.40$
withering	0	0	0.59

$\cos(\text{SaS}, \text{PaP}) = 0.6557 + 0.286 + 0 + 0$
 $= 0.9417$
 $= \underline{0.94}$

can assume written by same author. \therefore similarity is highest.

$\cos(\text{SaS}, \text{WH}) = 0.4187 + 0.2444 + 0.136 + 0$
 $= 0.7991$
 $= \underline{0.80}$

$\cos(\text{PaP}, \text{WH}) = 0.4399 + 0.2585 + 0 + 0$
 $= 0.6984$
 $= \underline{0.70}$



26-08-24.

Q. Calc the cosine similarity of D_1, D_2, D_3 for the given term freq

term	D_1	D_2	D_3
College	100	57	12
Farewell	10	30	70
Election	50	20	80.

i) log freq weighting.

$$w = 1 + \log_{10} f$$

term	D_1	D_2	D_3
College	3	2.76	2.07
Farewell	2	2.48	2.85
Election	2.69	2.30	2.90
$\sqrt{\sum w^2}$	4.49	4.37	4.56

ii) length-normalized vector. $w_i / \sqrt{\sum w^2}$

term	D_1	D_2	D_3
College	0.67	0.63	0.45
Farewell	0.45	0.57	0.625
Election	0.59	0.53	0.64

26-08-24

$$\cos(D_1, D_2)$$

$$= 0.4221 + 0.2565 + 0.3127$$

$$= 0.9913$$

$$= \underline{0.99}$$

$$\cos(D_2, D_3)$$

$$= 0.2835 + 0.35625 + 0.3392$$

$$= 0.97895$$

$$= \underline{\underline{0.98}}$$

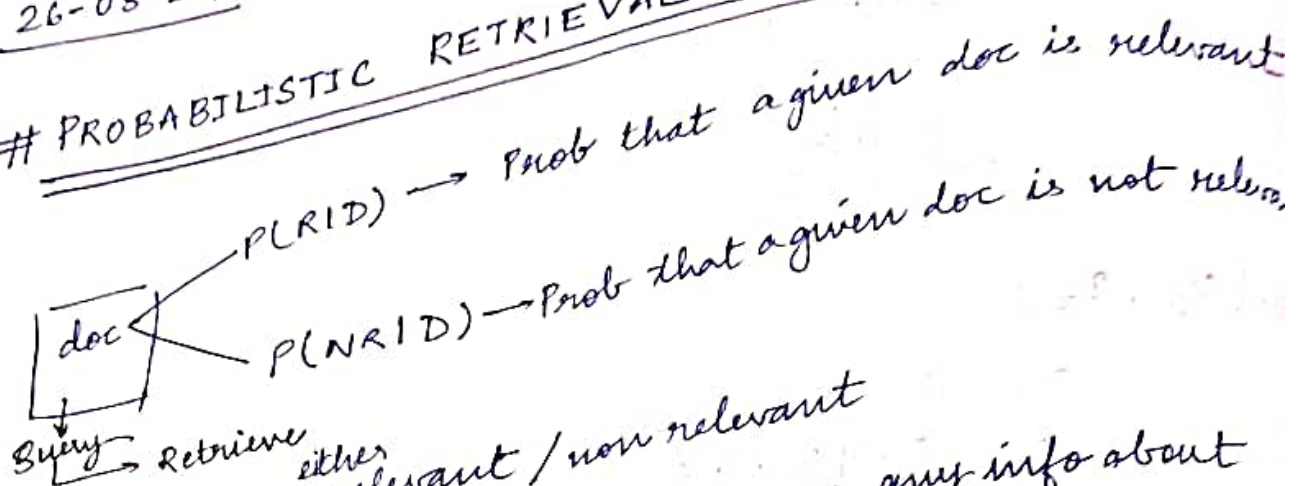
$$\cos(D_1, D_3)$$

$$= 0.3015 + 0.28125 + 0.3776$$

$$= \underline{0.96}$$

26-08-24

PROBABILISTIC RETRIEVAL MODEL



- i) A doc is ^{either} relevant / non relevant
- ii) Relevance of a doc does not convey any info about other doc being relevant.

Probabilistic Ranking principle (PRP)

- Ranking documents based on the decreasing probability of relevance to a query for the available data.
- Probability of relevance of a doc D for a given query, Q is

$P(R_Q = x | D)$

$x \in \{0, 1\}$

$1 \rightarrow$ relevant

$0 \rightarrow$ Not relevant

$$\frac{P(R_Q = 1 | D)}{P(R_Q = 0 | D)} \Rightarrow \text{Matching score of query \& document.}$$

$$= \frac{P(R_Q = 1) P(D | R_Q = 1)}{P(R_Q = 0) P(D | R_Q = 0)}$$

$P(D | R_Q)$
 \hookrightarrow prob. of retrieval of relevant doc

28-08-24

OKapi Retrieval Model

It is a sophisticated ranking func that builds on tradition TF-IDF method but emphasises the importance of less common word. It rank doc based on their relevant score, which is determined by follⁿ eqⁿ:

$$\sum_{i \in q} \log_e \frac{(q_i + 0.5)/(R - q_i + 0.5)}{(n_i - q_i + 0.5)/(N - n_i - R + q_i + 0.5)} \cdot \frac{(K_1 + 1)f_i}{K + f_i} \frac{(K_2 + 1)q_i}{K_2 + q_i}$$

where,

- q_i - no. of relevant doc containing item i
- n_i - no. of doc containing item i
- N - total no. of doc in collection
- R - no. of relevant doc for this query
- f_i - freq of item i in the doc.
- q_i - freq of item i in the query.

K_1, K_2, K - parametric parameter value set empirically

$$K = K_1 \left((1-b) + b \left(\frac{L_d}{L_{avg}} \right) \right)$$

\downarrow length of doc \downarrow avg. length of doc.

28-08-24

8) Query = "president" & "lincoln"
N = 500 000 doc.

$$k_1 = 1.2, b = 0.75$$

"president" occurs in 40,000 doc ($n_1 = 40\,000$)

"lincoln" occurs in 300 doc ($n_2 = 300$)

In a particular doc (D) that we are scoring, "president" occurs 15 times ($f_1 = 15$) & "lincoln" occurs 25 times ($f_2 = 25$).

The doc length is 90% of the average length ($L_d/L_{avg} = 0.9$)

$$K = k_1((1-b) + b(L_d/L_{avg}))$$

$b=0$; no normalization
 $b=1$; normalized

$$= 1.11$$

$$q_1 = q_2 = R = 0 \text{ (Assuming)}$$

$$n_1 = 40\,000$$

$$n_2 = 300$$

$$N = 500\,000$$

$$k_1 = 1.2$$

$$k_2 = 100$$

$$b = 0.75$$

$$K = 1.11$$

$$f_1 = 15$$

$$f_2 = 25$$

$$q_1 = 1$$

$$q_2 = 1$$

score

$$= \log_e \frac{(n_1 + 0.5) / (R - n_1 + 0.5)}{(n_1 - n_1 + 0.5) / (N - n_1 - R + n_1 + 0.5)} \left(\frac{(K_1 + 1) f_1}{K + f_1} \right) \frac{(k_2 + 1) 2 f_1}{K_2 + 2 f_1}$$

$$+ \log_e \frac{(n_2 + 0.5) / (R - n_2 + 0.5)}{(n_2 - n_2 + 0.5) / (N - n_2 - R + n_2 + 0.5)} \left(\frac{(K_1 + 1) f_2}{K + f_2} \right) \frac{k_2 + 1}{k_2 + 2 f_2}$$

$$= \log_e \frac{0.5 / 0.5}{40000.5 / 460000.5} \left(\frac{2.2 \times 15}{1.11 + 15} \right) \left(\frac{101}{101} \right)$$

$$+ \log_e \frac{0.5 / 0.5}{300.5 / 499700.5} \left(\frac{2.2 \times 25}{1.11 + 25} \right) \left(\frac{101}{101} \right)$$

$$= \left(\begin{matrix} 2.44 \\ 1.1016 \end{matrix} \times 2.05 \right)$$

$$+ \left(\begin{matrix} 7.42 \\ 3.2280 \end{matrix} \times 2.11 \right)$$

$$= 2.173 + 6.7942 = 5.002 + 15.66$$

$$= 8.9672$$

$$= 8.97$$

$$= 20.662$$

$$= 20.66$$

(3)



30-08-24

8. Consider:

query = "Information"

$N = 100$

"Information" word occurs in 37 doc
The doc length is 90% of avg length $(L_d/L_{avg}) = 0.9$
 $K_1 = 1.2$, $b = 0.75$, $K_2 = 100$, $K_3 = 1.11$

In a particular doc, D that we are scoring, we want to score,
"inf." occurs 12 times $= f_1 = 12$.

Calc score of doc using idf method.

Score =

$$\begin{aligned} & \log_e \frac{(n_i + 0.5) / (R - n_i + 0.5)}{(n_1 - n_i + 0.5) / (N - n_1 - R + n_i + 0.5)} \times \frac{(K_1 + 1) f_1}{K + f_1} \times \frac{(K_2 + 1)}{K_2 + 1} \\ &= \log_e \frac{63.5}{103.5} \times \frac{2.28 \times 12}{13.11} \times \frac{101}{101} \\ &= 0.5266 \times 2.013 \\ &= \underline{1.06} \end{aligned}$$

30-08-24

Performance Measure

IR system

TP \rightarrow doc retrieved is relevant
(true +ve)

TN \rightarrow doc not retrieved & it is not relevant
(true -ve)

FP \rightarrow doc retrieved & not relevant
as relevant
(false +ve)

FN \rightarrow Not retrieved & relevant
(relevant but not retrieved)
(false -ve)

	Relevant	Not Rel
Retrieved	TP	FP
Not Ret	FN	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision = fraction of retrieved doc that are relevant

$$P = \frac{TP}{TP + FP}$$

30-08-24

* Recall = fraction of relevant doc that is retrieved.

$$R = \frac{TP}{TP+FN}$$

* ~~#~~ f_1 measure
↓

Harmonic mean of P & R

$$f_1 = 2 \left(\frac{P \times R}{P+R} \right)$$

8.

	Rev	N/Rev
Ret	7	2
N/Ret	5	6

$$P = \frac{7}{7+2} = \frac{7}{9} = 0.777 = 0.78$$

$$R = \frac{7}{7+5} = \frac{7}{12} = 0.5833 = 0.583$$

$$Acc = \frac{7+6}{7+5+2+6} = \frac{13}{20} = 0.65$$

$$f_1 = 2 \left(\frac{7 \times 7}{7 \times 12} \right) = 2 \left(\frac{0.454}{1.363} \right) = 2 \times 0.333 = 0.666$$