

- 1) NO, it is not necessary for all zones to use the same Boolean Match function in weighted zone scoring, it can be beneficial to use different match functions for different zones based on their importance.
- 2) Assuming there are three zones and each zone has a boolean match function. The score for the document can be calculated as:-

$$\text{Score} = g_1 \times S_1 + g_2 \times S_2 + g_3 \times S_3$$

Here all distinct score values a document may get:

1) ~~0.2~~ $0.2 \times 0 + 0.31 \times 0 + 0.49 \times 0 = 0$

2) $0.2 \times 0 + 0.31 \times 0 + 1 \times 0.49 = 0.49$

3) $0.2 \times 0 + 0.31 \times 1 + 0.49 \times 0 = 0.31$

4) $0.2 \times 0 + 0.31 \times 1 + 0.49 \times 1 = 0.8$

5) $0.2 \times 1 + 0.31 \times 0 + 0.49 \times 0 = 0.2$

6) $0.2 \times 1 + 0.31 \times 0 + 0.49 \times 1 = 0.69$

7) $0.2 \times 1 + 0.31 \times 1 + 0.49 \times 0 = 0.51$

8) $0.2 \times 1 + 0.31 \times 1 + 0.49 \times 1 = 1$

3) ZONESCORE (List(ev))

float score[N] = [0]

constant g[1]

PZ-merge (List(ev))

while p (is not) NIL

Scores[docID(p)] = weightZone(p, g)

PZ-next(p)

return (scores)

4) Weight Zone(P_1, P_2, g)
 $S \leftarrow ()$
 $Scores[docID(P_1)] = 0$

for $i \leftarrow 1$ to 1

$S[i] \leftarrow BooleanScore(v, docID(P_1))$

$Scores[docID(P_1)] = Scores[docID(P_1)] + g[i] * S[i]$

5) By using the equation

$$g = \frac{n_{10r} \times n_{01n}}{n_{10r} + n_{10n} + n_{01r} + n_{01n}}$$

Given the sample training set

$n_{10r} = 3$ (from $\phi 1, \phi 3, \phi 5$)

$n_{10n} = 1$ (from $\phi 2$)

$n_{01r} = 0$

$n_{01n} = 1$ (from $\phi 4$)

So,

$$g = \frac{3 \times 1}{3 + 1 + 0 + 1} = 3/5 = 0.6$$

6) Phi g Relevance judgment

- 1 1R
- 2 3/4 NR
- 3 3/4 R
- 4 0NR
- 5 1 R
- 6 3/4 R
- 7 1/4 NR

7) In cases where $st(dt, art)$ and $sb(dt, art)$ have same values, score is independent of g and so these cases do not play a role in optimizing g .

8) idf always finite when
 $df_t \geq 1 \Rightarrow idf_t \leq \log N$.

9) $IDF(t) = \log(N/df(t))$

$IDF(t) = \log(N/N) \Rightarrow \log(1) = 0$

So, the IDF of a term that occurs in every document is 0. For a word that occurs in every document putting it in the stop list has the same effect as idf weighting: the word is ignored.

10)

	D.1	D.2	D.3
car	27	4	24
auto	3	33	0
insurance	0	33	29
best	14	0	17

the tf-idf weights for the terms:

	Doc1	Doc2	Doc3
car	44.55	6.6	39.6
auto	6.24	68.64	0
insurance	0	53.46	46.98
best	21	0	25.5

11) Yes, the tf-idf weight of a term in a document exceed 1.

2) For any base $b > 0$, $idf_t = \log_b(N/df_t)$
 $= (\log_b 10) * (\log_{10}(N/df_t))$
 $= c * (\log(N/df_t))$
 where c is constant.

$$tf-idf_{t,d,b} = tf_{t,d} * idf_t$$

$$= tf_{t,d} * c * (\log(N/df_t)) = c * tf-idf_{t,d}$$

$$Score(q, d, b) = \sum_w tf-idf_{t,b,d} = c * \sum_w tf-idf_{t,d}$$

So, changing the base changes the score by a factor $c = (\log_b 10)$

The relative scoring of documents remains unaffected by changing the base.

13) It is the number of bits in the boolean representation of the idf.

14) If "jealous" and "jealousy" both appear in a document, they would contribute to the TF of the stem "jealous".

• The modified TF for a stemmed term ts in a document d can be calculated as:

$$TF_{stem}(ts, d) = \frac{\text{Number of times stemmed term } ts \text{ appears in } d}{\text{Total number of term in } d}$$

• The modified IDF for a stemmed term ts can be calculated as:

$$IDF_{stem}(ts) = \log \left(\frac{N}{df_{stem}(ts)} \right)$$

15) Doc 1 = (0.897, 0.125, 0, 0.423)

Doc 2 = (0.076, 0.786, 0.613, 0)

Doc 3 = (0.595, 0, 0.706, 0.383)

16)

$$\text{Doc 1} = (0.897)^2 + (0.125)^2 + (0)^2 + (0.423)^2 = 0.999$$

$$\text{Doc 2} = (0.076)^2 + (0.786)^2 + (0.613)^2 + (0)^2 = 0.999$$

$$\text{Doc 3} = (0.595)^2 + (0)^2 + (0.706)^2 + (0.383)^2 = 0.999$$

Because they are normalized (unit) vectors.

$$18) \sum (q_i - w_i)^2 = \sum q_i^2 - 2 \sum q_i w_i + \sum w_i^2$$

$$\text{Thus,} \quad = 2(1 - \sum q_i w_i)$$

$$\sum (q_i - w_i)^2 < \sum (q_i - w_i)^2 \Leftrightarrow 2(1 - \sum q_i w_i) < 2(1 - \sum q_i w_i) < 2(1 - \sum q_i w_i) \Leftrightarrow \sum q_i w_i > \sum q_i w_i$$

19)

Word	query					document			
	tf	wf	df	idf	$q_i = wf \cdot idf$	tf	wf	$d_i = \frac{\text{normalize}}{wf}$	$q_i \cdot d_i$
digital	1	1	10,000	3	3	1	1	0.52	1.56
Video	0	0	100,000	2	0	1	1	0.52	0
Cameras	1	1	50,000	2.3	2.3	2	1.3	0.68	1.56

$$\text{Similarity, Score} = 1.56 + 1.56 = 3.12$$

$$\text{Normalize similarity score is also correct} = 3.12 / \text{length (query)}$$

$$= 3.12 / 3.78$$

$$= 0.825$$

20) $q_i = \text{affection}$

$$v(w) = (1, 0, 0), v(SoS) = (0.996, 0.087, 0.017)$$

$$v(PeP) = (0.993, 0.120, 0), v(WH) = (0.847, 0.466, 0.254)$$

$$v(w) \cdot v(SoS) = 0.996$$

$$v(w) \cdot v(PeP) = 0.993$$

$$v(w) \cdot v(WH) = 0.847$$

So, order of scores in this case is the reverse as in the query jealous gossip.

- 21) We can assign weights to query terms according to their idf in the collection, or use other Collection Statistics.
- 22) Omit this term from the query and proceed, its contribution to the dot product with any documents will be zero.
- 24) NTC weight contribution by coyotes insurance in any document vector = weight contribution by coyote + weight contribution by insurance = $0 + k$ (can be calculated) (since coyote does not occur in any document of the collection). The ctc contributed by coyote in the query vector need not to be calculated as the ntc weights for all documents is 0 for coyotes.