

TF \Rightarrow value over max in column

IDF $\Rightarrow \log_2 (\# \text{ docs} / \# \text{ docs word is in})$

$$w_{iq} = \left(0.5 + \frac{0.5 f_{iq}}{\max(f_{iq})} \right) \times \log \frac{N}{df_i}$$

Dot product and cosine similarity is a similarity check

$$\text{Query Cate} = \text{item TF \&F} \cdot \left(\frac{(\# \text{ times word in query})}{(\# \text{ words in query})} \right)$$

Dot product: $\text{sim}(\text{item in doc}, \text{query}) \rightarrow \text{sum}(\text{item} \cdot \text{query})$

Euclidean Distance: $\sqrt{\text{sum}((\text{doc item} - \text{query})^2)}$

Cosine Angles: $0-90^\circ$ map $0-1$, 1 means same
 0 totally different

Cosine Similarity: $\frac{\text{Sum of (doc item} \times \text{query)}}{\sqrt{\text{sum of square items in doc}} \sqrt{\text{sum of square items in query}}}$

$$\text{Precision} = \frac{\# \text{ Relevant items retrieved}}{\# \text{ retrieved items}} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{\# \text{ relevant items retrieved}}{\# \text{ relevant items}} = \frac{TP}{TP+FN}$$

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

Confusion matrix

What the system thinks

The truth

Total # of retrieved docs

↑ DCG changes ordering

	Relevant	Nonrelevant	Total
Retrieved	True positives (TP)	False positives (FP)	TP + FP
Not Retrieved	False negatives (FN)	True negatives (TN)	FN + TN
Total	TP + FN	FP + TN	TP + FP + FN + TN
	Total # of relevant docs	Total # of nonrelevant docs	Total # of docs in collection

F measure (F): the weighted harmonic mean of precision and recall, trading off precision versus recall.

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}, \text{ where } \beta^2 = \frac{1-\alpha}{\alpha}$$

where $\alpha \in [0, 1]$ and thus $\beta^2 \in [0, \infty]$.

The default balanced F measure equally weights precision and recall, which

means making $\alpha = 1/2$ or $\beta = 1$. $\rightarrow F_1$ measure $F_1 = \frac{2PR}{P+R}$

- $\beta < 1$: emphasize precision;
- $\beta > 1$: emphasize recall.

$$F_1@k = \frac{2 \times (\text{Precision}@k) \times (\text{Recall}@k)}{(\text{Precision}@k) + (\text{Recall}@k)}$$

$$\text{Recall}@k = \frac{\# \text{ relevant items @ } k}{\text{total \# relevant items}} = \frac{TP@k}{TP + FN}$$

All @k

$$\text{Precision @ } k = \frac{\# \text{ relevant items @ } k}{\# \text{ retrieved items @ } k} = \frac{TP}{TP + FN}$$

Mean reciprocal rank is defined as: $MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{rank}_i}$

- $|Q|$: the total number of queries
- rank_i : the rank of the first relevant result retrieved in query i .

$$AP = \frac{\text{sum (Precision @ } k \text{ when relevant)}}{\text{total \# relevant items}}$$

$$MAP = \frac{1}{|Q|} \sum_{q=1}^{|Q|} AP(q)$$

Q : total # of queries

$AP(q)$: average precision for query q .

\rightarrow only 1 result

Cumulative gain @ k
is sum of
relevance until
@ k

DCG @ k is sum of
(relevance score / $\log_2(k+1)$)

$$NDCG = DCG / IDCG$$