TF is feq of word / max word freq in doc, idf is $\log \dfrac{\text{num doc total}}{\text{num docs in}}$

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

Eluciden Distance: $\sqrt{\text{sum}((\text{doc item} - \text{query})^2)} \rightarrow$ large diff length

Cosin Angles: $0 - 90°$ map $0 - 1$, 1 means same
$0$ totally different

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

Precision $= \dfrac{\text{\# Relavant items retrieved}}{\text{\# retrieved items}} = \dfrac{TP}{TP+FP}$

Recall $= \dfrac{\text{\# relavant items retrived}}{\text{total \# relavant items)}} = \dfrac{TP}{TP+FN}$

Accuracy $= \dfrac{TP + TN}{TP+FP+FN+TN} = \dfrac{\text{\# true}}{\text{total}}$

Query $= (\text{num times item appears / max times item appears}) \times \text{tfidf}$

| | Relavant | Nonrelavant | total | |
|---|---|---|---|---|
| retrived | TP | FP | TP+FP | → total # retrived |
| Not retrived | FN | TN | FN+TN | |
| total | TP+FN | FP+TN | total ↕ | |
| | ↑ | ↑ | ↑ | |
| | \#relavant) | \# nonrelavt | total \# | |

$$F_{measure} = \frac{1}{\alpha \frac{1}{P} + (1-\alpha)\frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

Reciprocal Rank $= 1/$ pos relevant

$$F_1 = \frac{2PR}{P+R} \rightarrow \text{this is when } \alpha = \frac{1}{2} \text{ or } \beta = 1$$

$$\text{Precision @} k = \frac{TP@k}{TP@k + FP@k}$$

Average Precision

Precision @ relevant
#of relevant total

$$\text{Recall @} k = \frac{TP@k}{TP + FN}$$

Initial Rank $1/\#$ items

New page rank is $\frac{PR \text{ of input}}{num \ outputs \ of \ item}$ Sum

$$MMR = \frac{1}{total \ \# \ query} \left( \sum_{i=1}^{\#query} \frac{1}{rank_i} \right)$$

If a sink you distribute page ranks equally on all other items, after new calc think iteration $i$ and $i'$

Robots.txt specifies rules of a webpage → get by /robots.txt

VSM assumes independence, needs to be precise

Crawl BFS, DFS

MAP is average precision on multiple queries

One pagerank can't tell how important a page is

web crawler like librarian → places items index for later

Freshness is accuracy of copy, Age is how outdated a copy is

Uniform policy → revisit all pages with the same frequency

Proportional Policy → revisit pages with more changes more frequently

Web crawling is search, web scraping is extracting data

avoid duplication

+create index for search | Comes from is http://www.os2.com/os2/