termID → documentID
(docID)

reducing the dimension
→ implies loss of info

Compression
→ loosy
→ loseless

(Esperanto)

Heaps' Law

rate of growth

## Jaccard Coefficient

A = This (is) (a) (test)³
B = (A) (test) (is) conducted.   $J(A,B) = \dfrac{A \cap B}{A \cup B} = \dfrac{3}{5}$

## Term frequency
(Document freq.)

~~Log weight (frequenc~~

## Log freq. weighting

$$w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d} & , \quad \text{if } tf_{t,d} > 0 \\ 0 & , \quad \text{otherwise} \end{cases}$$

log freq.
wt. in
document d

D = An apple a day keep you away from dacto
Apple is good for health. Nit are pr in
Apple.

(apple)

$$= 1 + \log_{10} 3$$
$$= 1 + 0.4771$$
$$= 1.4771$$

Log ~~term~~ freq- weighting

$$\text{Score} = \sum_{t \in (q \cap d)} \left(1 + \log_{10} tf_{t,d}\right)$$

IDF $\longrightarrow$ inverse doc. freq.

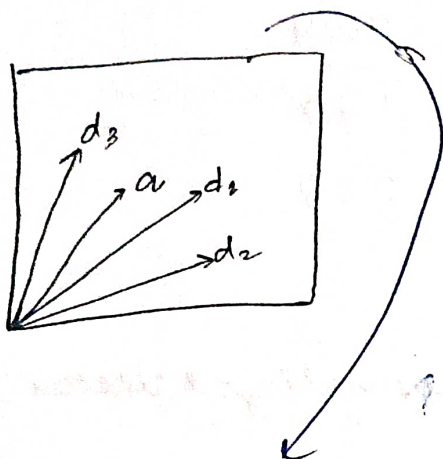$$idf_t = \log_{10}\left(\frac{N}{df_t}\right)$$  tf   TF- IDF

IDF weight of a term t in doc d is

$$W_{t,d} = \log\left(1 + tf_{t,d}\right) \times \log_{10}\left(\frac{N}{df_t}\right)$$

* 7 Score of a doc. for a given query —

$$\text{Score}(q,d) = \sum_{t \in q \cap d} tf \cdot idf_{t,d}$$

## Document as vector



Term as diff. axes

similarity.

1) distance

distance $\frac{1}{\alpha}$ similarity

less dist. more simila[r]

Cosine similarity.

cosine similarity bet^n 3 docs.

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

PaP → pride & prejudice
SaS → sense & sensibility
WH → wuthering height

ans

# Log frequency weighting

| term | SaS | PaP | WH |
|------|-----|-----|-----|
| affection | 3.06 | 2.76 | 2.3 |
| jealous | 2 | 1.04 | 2.04 |
| gossip | 1.3 | ⓪ | 1.77 |
| wuthering | ① | ⓪ | 2.57 |

cosine similarity for length normalized vector.

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d}$$

$$= \sum_{i=1}^{|V|} q_i d_i$$

## Length Normalization

$$\frac{3.06}{\sqrt{(3.06)^2 + 2^2 + (1.30)^2}}$$

$$\frac{2.76}{\sqrt{(2.76)^2 + (1.04)^2}}$$

| term | SAS | PAP | WH |
|------|-----|-----|-----|
| affection | 0.79 | 0.83 | 0.524 |
| jealous | 0.515 | 0.55 | 0.465 |
| gossip | 0.835 | 0 | 0.405 |
| wuthering | 0 | 0 | 0.588 |

$\cos(SAS, PAP)$

$\Rightarrow \quad 0.79 \times 0.83 + 0.515 \times 0.55 + 0.7 + 0$

$= \quad 0.93$

$\cos(PAP, WH) = 0.83 \times 0.524 + 0.55 \times 0.465 + 0 + 0$
$\qquad = \quad 0.6697$

$\cos(SAS, WH) = \quad 0.79 \times 0.524 + 0.515 \times 0.468 + 0.335 \times 0.405 + 0$

$\qquad = \quad 0.786$

| item | D1 | D2 | D3 |
|------|-----|-----|-----|
| college | 100 | 57 | 12 |
| Election | 50 | 80 | 70 |
| farewell | 10 | 20 | 80 |

Calculate cosine similarity of $\cos(D_1, D_2) = ?$
$\cos(D_1, D_3) = ?$
$\cos(D_2, D_3) = ?$

$$\cos(D_1, D_2) = 100 \times 57 + 50 \times 30 + 10 \times 20$$
$$= 5700 + 1500 + 200$$
$$= 7400$$

$$\cos(D_1, D_3) = 100 \times 12 + 50 \times 70 + 10 \times 80$$
$$= 1200 + 3500 + 800$$
$$= 5500$$

$$\cos(D_2, D_3) = 57 \times 12 + 30 \times 70 + 20 \times 80$$
$$= 684 + 2100 + 1600$$

57
114
57
684

3780
684

$$\Rightarrow \frac{2.07}{\sqrt{(2.07)^2 + (2.04)^2 + (2.9)^2}}$$
$$= \frac{2.84}{10.57}$$
$$= \frac{2.9}{18.54}$$

$$\Rightarrow \frac{2.07}{10.54}$$

$$\Rightarrow \frac{3}{\sqrt{3^2 + (2.69)^2 + 2^2}}$$

$$\Rightarrow \boxed{\frac{3}{14}}$$

$$= \frac{2.69}{14}$$

$$\Rightarrow \frac{2.69}{\sqrt{3^2 + (2.69)^2 + 2^2}}$$

$$\Rightarrow \frac{2}{14}$$

$$\Rightarrow \frac{2}{\sqrt{3^2 + (2.69)^2 + 2^2}}$$

$$\Rightarrow \frac{2.78}{\sqrt{(2.78)^2 + (2.47)^2 + (2.3)^2}}$$

$$\Rightarrow \frac{2.47}{}$$

$$\Rightarrow \frac{2.3}{}$$

$$= \frac{2.78}{14.17}$$

log freq. weighting

↓

Length normalization

~~truncated~~

Cosine

log freq. weight -

| item | $D_1$ | $D_2$ | $D_e$ |
|---|---|---|---|
| College | 3 | 2.75 | 2.07 |
| Election | 2.69 | 2.47 | 2.04 |
| farewell | 2 | 2.30 | 2.90 |

Length normalization

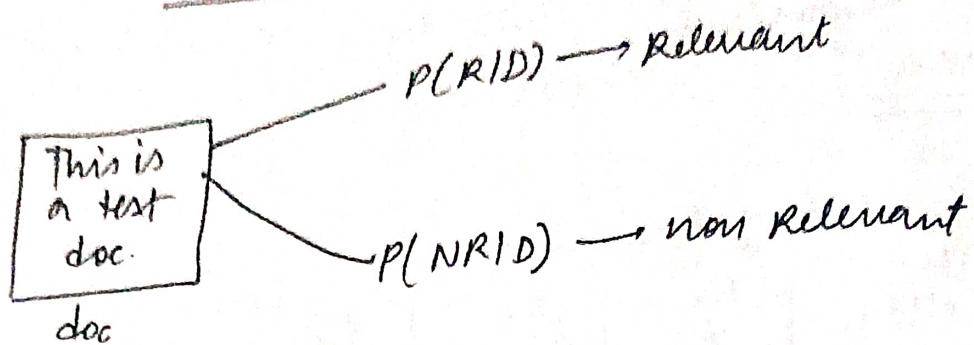| item | $P_1$ | $D_2$ | $D_3$ |
|---|---|---|---|
| College | 0.2 | 0.19 | 0.11 |
| Election | 0.19 | 0.17 | 0.15 |
| farewell | 0.14 | 0.16 | 0.16 |

Cosine

$\cos(D_1, D_2) = 0.2 \times 0.19 + 0.19 \times 0.17 + 0.14 \times 0.16$
$= 0.0927$

$\cos(D_1, D_3) = 0.2 \times 0.11 + 0.19 \times 0.15 + 0.14 \times 0.16$
$= 0.0729$

$\cos(D_2, D_3) = 0.19 \times 0.11 + 0.17 \times 0.15 + 0.16 \times 0.16$
$= 0.072$

$D_1$ & $D_2$ has maximum relevance.

# Probabilistic Information Retrieval System (Model)



```
              ┌─────── P(R|D) ──→ Relevant
  ┌────────┐ ─┤
  │This is │  └─────── P(NR|D) ──→ non Relevant
  │a test  │
  │doc.    │
  └────────┘
    doc
```

→ examine the relevance of the document.

→ probabilistic Ranking principle.

→ decreasing probability of relevance.

→ highest probability at the top.

## Probabilistic Ranking principle

N documents

### probability of Relevance -

$$P(R = X/D)$$

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

$X = 1 \longrightarrow$ Relevant

$X = 0 \longrightarrow$ Non Relevant

$P(R = 1 | D)$

$P(R = 0 | D)$

The Score of Matching

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

$$\boxed{P(R|D) = \frac{P(D/R) \cdot P(R)}{P(D)}}$$

# Okapi Model

→ probabilistic IR Model

$$TDF \longrightarrow \log_{10}\left(\frac{N}{df_t}\right)$$

$$\boxed{\log_e\left(\frac{N}{df_b}\right) \cdot \frac{(K_1 + 1)\, tf_{t,d}}{K_1\left[(1-b) + b\left(\frac{L_d}{L_{avg.}}\right)\right] + tf_{t,d}}}$$

$tf \longrightarrow$ term frequency

$L_d \longrightarrow$ Length of the doc.

$L_{avg} \longrightarrow$ avg. length of the doc.

$K_1 \longrightarrow$ Tuning parameter controlling the doc. term frequency scaling

$b \longrightarrow$ Tuning para for doc. length

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

$$\frac{(K_2 + 1)\, qf_i}{K_2 + qf_i}$$

$r_i \longrightarrow$ relevance of $i^{th}$ term

**9.** Query: "President" & "Lincoln"

$R = r = 0$

$N = 5000,000$ doc

$qf_i (i=1) = $ President occurs in query $= 1$

"president" word occurs in 40000 doc.

"lincoln" word occurs in 300 doc

$qf_i (i=2) = $ Lincoln occurs in query $= 1$

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 aug. length. i.e $\dfrac{Ld}{Laug} = 0.9$

$K_1 = 1.2$  $b = 0.75$  &  $K_2 = 100$

$$\boxed{K = K_1(1-b) + b\left(\dfrac{Ld}{Laug}\right)}$$

$$K = 1.11$$

$$\Rightarrow \log_e \dfrac{(0+0.5)/(0-0+0.5)}{(40000-0+0.5)/(5000000-40000-0+0+0.5)} \cdot \dfrac{(1.2+1)15}{1.11+15} \cdot \dfrac{(100+1)\times 1}{(100+1)}$$

$$+ \log_e \dfrac{(0+0.5)/(0-0+0.5)}{(300-0+0.5)/(5000000-300-0+0+0.5)} \cdot \dfrac{(1.2+1)\times 25}{1.11+25} \cdot \dfrac{(100+1)\times 1}{100+1}$$

Now ~~// // //~~

$$\Rightarrow \log_e \dfrac{1}{0.008} \cdot \dfrac{2.2\times 15}{16.11} \times \dfrac{101'}{101'} + \log_e \dfrac{1}{0.00006} \times \dfrac{2.2\times 25}{26.11} \times \dfrac{101'}{101'}$$

$$\Rightarrow \log_e \dfrac{1000}{8}^{125} \cdot 2.04 + \log_e \dfrac{100000}{6} \times 2.10$$

$$\Rightarrow \log_e 125 \times 2.04 + 9.72 \times 2.106$$

$$= 4.82 \times 2.04 + 20.47$$

$$- 9.04 + 20.47$$

$$= 30.350$$

# Language Modeling

→ If query is more refine.

→ Simplest model is **unigram** model.

   n-gram

**5 words** → In a document collection

→ (0.2, 0.1, 0.35, 0.25, 0.1)

"cat", "rain", "dog", "jump", "the"

Probability of occurring 2 words.

"cat rain" = 0.2×0.1 = 0.02

"cat jump" = 0.2×0.25 = 0.5

→ $qf_t = 1$

query = 'Apple' ⇒ q

N = 100

Apple word occurs in 37 doc

In a particular document (D) that we want to score.
apple occurs 12 times.

The doc-length is 90% of the average length

$\left( \frac{L_d}{L_{avg.}} \right) = 0.9$

$K_1 = 1.2$, $b = 0.75$

$0.25 + 0.75 × 0.9$

$K_2 = 100$

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

$= 1.11$

cal. score (q, D) using Okapi Method -

$$\log_e \left( \frac{N}{df_t} \right) \cdot \frac{(k_1 + 1)\, t_{f_{t,d}}}{K_1 \left[ (1-b) + b \left( \frac{L_d}{L_{avg}} \right) \right] + t_{f_{t,d}}}$$

$$\Rightarrow \log_e \left( \frac{100}{37} \right) \cdot \frac{(1.2 + 1) × 12}{1.11 + 12} \Rightarrow 0.99 × \frac{2.2 × 12}{13.11} \Rightarrow 0.99 × 2.01$$

$$= 1.9935$$

$$\log_e \frac{(0+0.5)/(0-0+0.5)}{(37-0+0.5)/(100-37-0+0+0.5)} \times \frac{(1.2+1)\,12}{1.11+12} \times \frac{(100+1)\,1}{(100+1)}$$

$$\Rightarrow \log_e \frac{63.5}{37.5} \times \frac{26.4}{18.11} \times 1$$

$$= \log_e(1.69)$$

$$= 0.52 \times 2.01$$

$$\Rightarrow 1.05$$

## Performance measure

|  | Relevant | N.R |
|---|---|---|
| Retrieved | TP | FP |
| Not Retrieved | FN | TN |

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

$P \longrightarrow$ fraction of retrieved document that are actually relevant.

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

$R \longrightarrow$ fraction of relevant document

|  | Rev. | NRev |
|---|---|---|
| Ret | 5 | 3 |
| NR | 2 | 5 |

$$Acc = \frac{10}{15}$$

$$= \frac{2}{3} = 66.67\%$$

$P = \#$ Rev. item $/\#$ Ret. $= \frac{5}{8} = 0.62\,8$

$R = \#$ Rev item retrieved $/$ Rev $= \frac{5}{7} = 0.74$

q - query
relevant doc = 20

relevant retrieved =

$$\frac{15\,doc}{5\,\text{(N.R)}}$$ N.Rel

Harmonic mean of P & R is known as

$f_1 \rightarrow$ measure

$$F1 = 2\,\frac{P * R}{P + R}$$