

IR: static (fixed-~~data~~ information range)

IF: on-going (stream).

(can add ~~some~~ words)

pre-processing:
(NLP).

(P) set

Use same
pre-processing

Q (set).

(underlying) docu.

(key words, phrases ...)

comparison

paragram → bag of words.
represent to

result.

Stemming:

dis: over stemming

sking

lokies.

absorb
absorp

partial string
match

abcdef
abcdefg

stop-words:

dis: someone has meaning in phrase.

"The the" someone has certain meaning.

Thesaurus construction:

thesis

↔ dissertation

or thesis

related words:

messi

{ soccer
barcelona
psg.

IR.

D. (pseudo-see).

(could be a bag of words)

(set of terms stems)

Q

framework (mathematical)

Ranking $\leftarrow \text{sim}(\text{query}, \text{doc})$

indexing

$t_1 \ t_2 \ \dots \ t_n$

$d_1 \ w_{11} \ w_{21}$

doc 2

d_3

w_{11} : how well the term
describe this doc.
relative to.

↓

how to assign weights.

this

matrix could be huge / sparse \Rightarrow it could be inefficient.

(weights can be score in vectors, some of them could be 0).

how to ~~see~~ find a good weight for each term?

Boolean Model.

binary matrix

	t_1	t_2	t_n
d_1	0	0	1
d_2			
d_3			
\vdots			

t set

$$t_1 \cap t_2$$

$$t_1 \cup t_2$$

disadv.: no difference in weights.

coffee

d_1

take no account to d_2 .

frequency

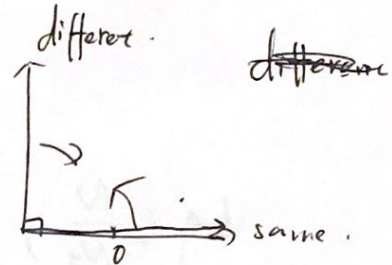
each term vectors are independent

Vector space Model.

	t_1	t_2	\dots	t_n
\vec{d} d_1				
d_2				
d_3				
\vdots				
\vec{q}	0.1	0.6	\dots	0.8

$$0 \leq w_{ij} \leq 1. \quad (\text{default}).$$

Similarity



$$\text{sim}(q, d) = \frac{\sum}{\sum}$$

calculate instance

\sum
already
calculated.

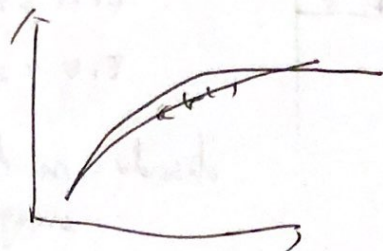
$$\sum_{t \in q \cap d}$$

return the short documents which contain my query.

normalise based on the length of query and documents.

tf : local

idf : global.



function of frequency.

1 → 2 difference is huge

50 → 51

less,

$f(\text{freq})$

tf x idf
(local) (global).

normalisation

→ e.g. if coffee appears in all doc.

f_{ij}

$\log\left(\frac{N}{N_{ij}}\right)$

↓
millions of
terms
frequency

↓
idf

(millions of documents).