Stop lists.

# Motivation & Methods of Stemming

Calculating TF-IDF Simularity

- Stop words Removal : Remove 'Noise words' from text
  
  find examples
  
  Contribute no info. to info. Retreival Process

- Stemming : Remove irrelavent differences.
  
  Different forms of the same word.
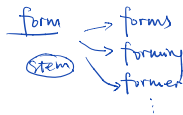
- Semantic Relationships.

Stemming     • A query. & a doc. contain different forms of the same word

(Morphology)     Related.

- Remove surface markings.
  Reveal basic form.

form    → forms
        → forming
Stem    → former
        :

- Replace. words with its eq. classes of words
                    Stems
  
  ⇒ Reduce number of different words
     Increase the number of instances of each token.

- Question raised : Not all words obey regular rules
  - Solution : Identify sub-pattern of letters
           ⇒ devise rules to deal with patterns

Stemmer : implements stemming algorithm.
    Using stemer ⇒ Reduce vocab. size   10% ~ 50%

Stop words     : vital for grammar
    "Noise"     useless for idenditifying the content

- Specified in a text file   stop list

Matching    Query q.    Document. d.

    Sim (q, d)  ⇒ def. Simularity
    • No. terms common to q and d.
      how
    • useful. is to common term
        e.g. "the" and. "magnesium"

## IDF weighting : Measuring Significance

inverse document frequency

$$IDF(t) = \log\left(\frac{ND}{ND_t}\right)$$

total num. of Docs.

# Docs include t.

- case ① t occurs in every docs      $ND = ND_t$.
  $\log(1) = 0$

② t in a few docs      $ND > ND_t$.
  $\log\left(\frac{ND}{ND_t}\right) > 0$.

* ignore the occurrence frequency within each doc.

Document length.


## TF - IDF weight
term frequency - Inverse Document frequency

$$\underline{W_{td}} = f_{td} \cdot IDF(t)$$
(Number of times t occurs in d) · Inverse frequency of t

doc.

$f_{td}$ Large: often occur in d.
$IDF(t)$ large: occur in few docs


## Query weights
- Long query $\underline{q}$   treat as document
  $$W_{tq} = f_{tq} \cdot IDF(t)$$
  # times t in q

- Short query $\underline{q}$
  $$W_{tq} = IDF(t)$$


## TF -IDF Similarity   between q and d.
$$Sim(q,d) = \frac{\sum_{t\, in\, q\, and\, d} W_{td} \cdot W_{tq}}{\|d\| \cdot \|q\|}$$

For every t in both q and d.
- Calculate document TF-IDF weight   ⎫ Sum the product
  Calculate query weight

# Document length

$$\text{Len}(d) = \|d\| = \sqrt{\sum_{t \in d} w_{td}^2}$$
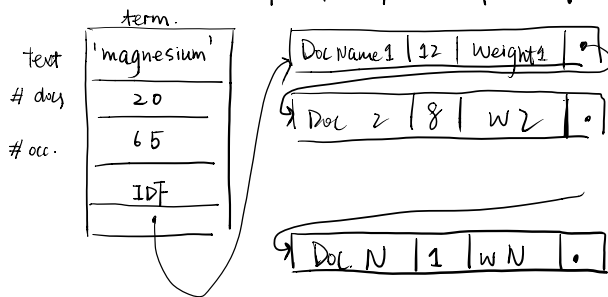
$f_{td} \cdot IDF(t)$

Vector space
$$\|X\| = \sqrt{X_1^2 + X_2^2 + X_3^2}$$

# Document Index
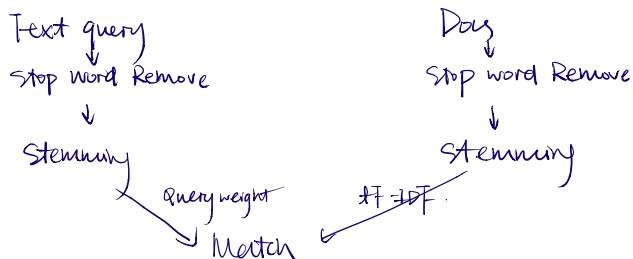
Document Index: Speed up computation of $\text{Sim}(q, d.)$

|                | term.         |
|----------------|---------------|
| text           | 'magnesium'   |
| # docs         | 20            |
| # occ.         | 65            |
|                | IDF           |

Doc Name 1 | 12 | Weight 1 | •

Doc 2 | 8 | W 2 | •

Doc N | 1 | w N | •

?

~~Docs are ordered by the occurence of t~~

"unique"

- Practical:
- terms: ordered in decreasing IDF

- For each term: Docs : decreasing weight

"if for t, its often and unique" ⇒ weights higher.

- For each term in Query
  - identify ~ in index.         ⇒ 相加
  - Increment similarity scores
  - Stop when weights fall below threshold.

*(while calculating Similarity)*

# IR Process

Text query
↓
Stop word Remove
↓
Stemming
↓
Query weight
↓
Match

Doc
↓
Stop word Remove
↓
Stemming
↓
打 IDF.

# Summary of Points

- Stop Lists
- Stemming
- TF-IDF weight for Docs
- Query weight
- "Lenghth" of doc. Len.(d)
- TF-IDF similarity

- Document index ⇒ speed up Sim. calculation.