

The term vocabulary and postings lists

lightsilver

Nov 9th

Obtain Character Sequence

- Convert byte sequence into a linear sequence of characters.
 - Complex to determine the encoding
 - Machine learning
 - Heuristic
 - User select
 - Using metadata
 - Complex according to file format
 - Markup in XML
 - Licensing software library

Document Unit

- File interface
 - Spilt attachments from mails
 - Combine separate HTML pages into single PPT
- Granularity level
 - Large: Get spurious matches, rare relevant result
 - Small: Miss important passage between documents
 - Precision / Recall tradeoff

Tokenization

Input: Friends, Romans, Countrymen, lend me your ears;

Output:

Friends	Romans	Countrymen	lend	me	your	ears
---------	--------	------------	------	----	------	------

- Simply spilt by white space and punctuation?
 - Larger vocabulary
 - Language
 - ...

Who is most likely correct 1.0

- O'Neil
 - neill
 - oneill
 - o'neil
 - o' & neil
 - o & neil
- Aren't
 - aren't
 - arent
 - are & n't
 - aren & t
- Type
 - class of all tokens containing the same character sequence.

Typed issue

- What if do the exact same tokenization of document and query words, generally by processing queries with the same tokenizer.

Typed issue

- What if do the exact same tokenization of document and query words, by processing queries with the same tokenizer.
- Always matches using postings lists?

Who is most likely correct 2.0

- C#
- bgisgood@gmail.com
- 10.13.122.223
- 1Z9999W99845399981
- <http://mail.mstczju.org>
- Semantic partition
 - Expend vocabulary / Restrict query input tradeoff
 - Metadata

Who is most likely correct 3.0

- Lebensversicherungsgesellschaftsangestellter
 - subdivided into multiple words
- 莎拉波娃现在居住在美国东南部
 - --
- 和尚
 - Monk or and + still?
- Language identification
 - word segmentation

Stop words

- Of little value in helping select documents matching a user need are excluded from the vocabulary entirely.
- Determined by collection frequency.
 - > Stop list

a	an	and	are	as	at	be	by	for	from
has	he	in	is	it	its	of	on	that	the
to	was	were	will	with					

Who is most likely correct 3.0

- “President of the United States”
 - “President” AND “United States”
- “To be or not to be”
 - Empty query. -_-
- Trend
 - Smaller stop list(precise > efficiency)

Normalization

- Canonicalize tokens despite superficial differences
- Equivalence Class
 - remove characters
 - anti-discriminatory and antidiscriminatory
 - maintain relations between unnormalized tokens
 - car and automobile
 - Index expansion(space)
 - Query disjunction(time)
 - pluskid / p1uskid

Who is most likely correct 4.0

- Query term Terms in documents that should be matched
 - Windows Windows
 - windows Windows, windows, window
 - window window, windows
- USA & U.S.A
 - ok
- C.A.T & cat
 - =_ =
- naïve & naïve
- 3/12/91
 - US Mar. 12, 1991, whereas in Europe it is 3 Dec 1991.

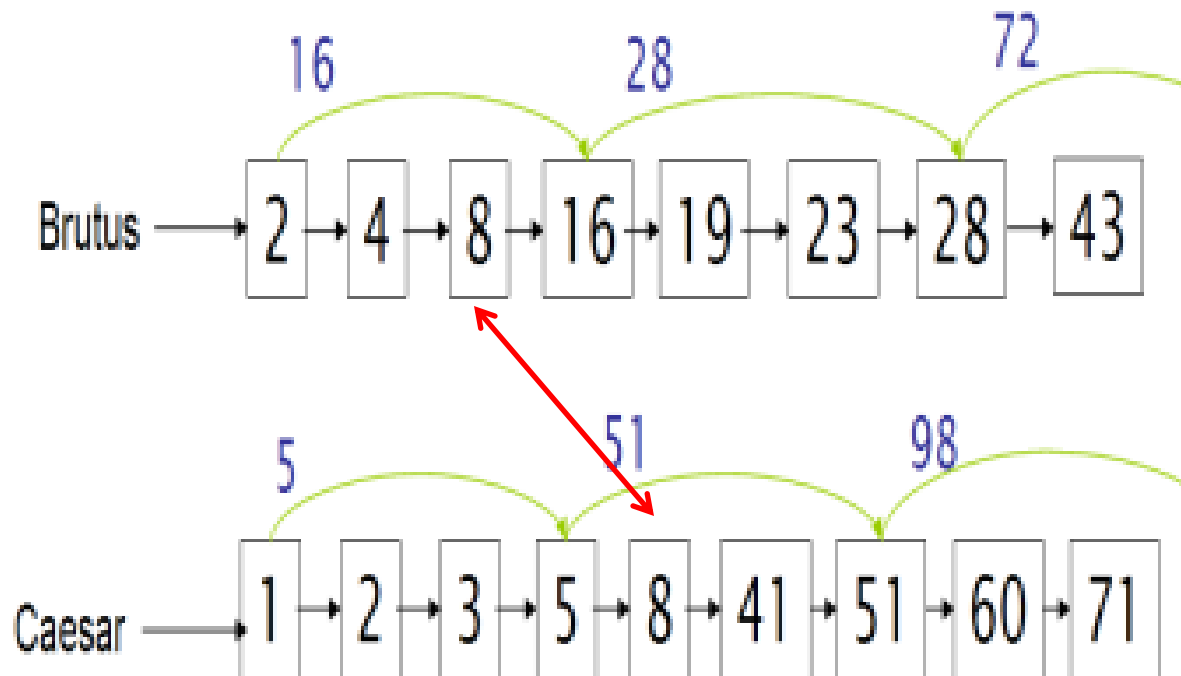
Semming&Limmatization

- Reduce inflectional forms and derivationally related forms of a word to a common base form.
- the boy's cars are different colors \Rightarrow
- the boy car be differ color
- Plug-in processing index.
- Porter's algorithm \rightarrow Semming

Who is most likely correct 5.0

- Stemming increases recall while harming precision.
- Porter semmer:
 - operate operating operates operation operative operatives operational
→ oper
- Lemmatizer:
 - operating system → operate system

Skip Pointer



- Only for original postings lists.
- Only for “and” queries.

Skip Pointer cont.

- Where to place skip pointers? → Tradeoff
 - More skips
 - Shorter spans
 - More comparisons
 - Space occupancy
 - $P^{(1/2)}$ evenly-spaced skip pointers from practice
- What about updates?

Fewer skips

Fewer pointer comparisons

Longer spans

Fewer opportunity to skip

Phrase Query

- Double quotes syntax (“stanford university”) doesn’t match this.
 - The inventor Stanford Ovshinsky never went to university.
- Proximity Weight Ranked Retrieval ?
- Biword index
- Positional index

Biword Index

- “stanford university palo alto” →
 - “stanford university” & “university palo” & “palo alto”
- “renegotiation of the constitution”
- → $N \times N$
- → NN
- → “renegotiation constitution”

Who is most likely correct 6.0

- “cost overruns on a power plant”
 - “cost overruns”
 - “overruns power” → omit this?
 - “power plant”
- Long phrases storage size
- Less recall

Positional Index

- Frequency for weighting to, 993427:
 - 1, 6: {7, 18, 33, 72, 86, 231};
 - 2, 5: {1, 17, 74, 222, 255};
 - 4, 5: {8, 16, 190, 429, 433};
 - 5, 2: {363, 367};
 - 7, 3: {13, 23, 191}; ...
- Further restrict the list of possible candidates
- Check compatible positions of appearance
- “to be or not to be”

to: {...; 4: {..., 429, 433}; ...}
 be: {...; 4: {..., 430, 434}; ...}

Algorithm

POSITIONALINTERSECT(p_1, p_2, k)

```
1  answer  $\leftarrow \langle \rangle$ 
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$ 
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$ 
4      then  $l \leftarrow \langle \rangle$ 
5           $pp_1 \leftarrow \text{positions}(p_1)$ 
6           $pp_2 \leftarrow \text{positions}(p_2)$ 
7          while  $pp_1 \neq \text{NIL}$ 
8          do while  $pp_2 \neq \text{NIL}$ 
9              do if  $|\text{pos}(pp_1) - \text{pos}(pp_2)| \leq k$ 
10                 then  $\text{ADD}(l, \text{pos}(pp_2))$ 
11                 else if  $\text{pos}(pp_2) > \text{pos}(pp_1)$ 
12                     then break
13                      $pp_2 \leftarrow \text{next}(pp_2)$ 
14                 while  $l \neq \langle \rangle$  and  $|l[0] - \text{pos}(pp_1)| > k$ 
15                     do  $\text{DELETE}(l[0])$ 
16                     for each  $ps \in l$ 
17                         do  $\text{ADD}(\text{answer}, (\text{docID}(p_1), \text{pos}(pp_1), ps))$ 
18                      $pp_1 \leftarrow \text{next}(pp_1)$ 
19              $p_1 \leftarrow \text{next}(p_1)$ 
20              $p_2 \leftarrow \text{next}(p_2)$ 
21         else if  $\text{docID}(p_1) < \text{docID}(p_2)$ 
22             then  $p_1 \leftarrow \text{next}(p_1)$ 
23             else  $p_2 \leftarrow \text{next}(p_2)$ 
24 return answer
```

- Position w/o compression
- Complexity
- BI PI combined
 - “Michael Jackson”
 - “Britney Spears”
 - “The Who”
- Improvement: Next Word Index