

Text Preprocessing

Reference: Introduction to Information Retrieval
by C. Manning, P. Raghavan, H. Schütze

Parsing a document

- What format is it in?
 - pdf/word/excel/html?
- What language is it in?
- What character set is in use?
 - (CP1252, UTF-8, ...)

These tasks are often done heuristically ...

Complications: Format/language

- Documents being indexed can include docs from many different languages
 - A single index may contain terms from many languages.
- Sometimes a document or its components can contain multiple languages/formats
 - Chinese email with an English pdf attachment.
 - French email quote clauses from an English-language contract
- There are commercial and open source libraries that can handle a lot of this stuff

Tokenization

- Input: “*Friends, Romans and Countrymen*”
- Output: Tokens
 - *Friends*
 - *Romans*
 - *Countrymen*
- A **token** is an instance of a sequence of characters
- Each such token is now a *candidate* for an index entry, after further processing
 - Described below
- But what are valid tokens to emit?

Tokenization

- Issues in tokenization:
 - *Finland's capital* →
Finland AND *s*? *Finlands*? *Finland's*?
 - *Hewlett-Packard* → *Hewlett* and *Packard* as two tokens?
 - *typical solution*: break up hyphenated sequence.
 - *co-education*
 - *lowercase, lower-case, lower case* ?
 - It can be effective to get the user to put in possible hyphens
 - *San Francisco*: one token or two?
 - How do you decide it is one token?

Numbers

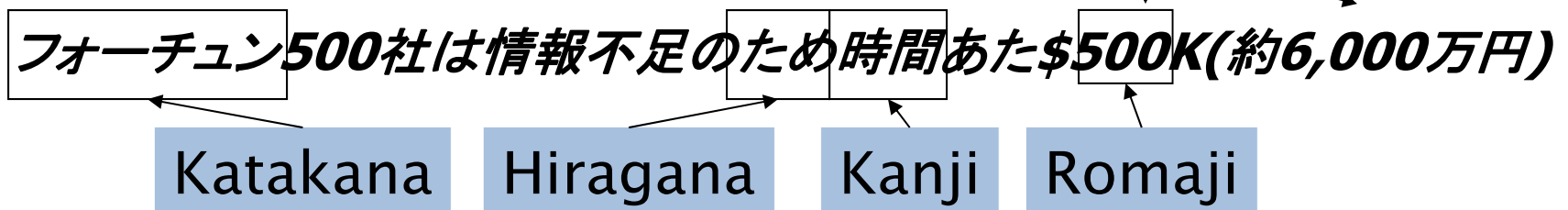
- ***3/20/91***
- ***55 B.C.***
- ***B-52***
- ***My PGP key is 324a3df234cb23e***
- ***(800) 234-2333***
 - Older IR systems do not index numbers
 - But often very useful: think about things like looking up error codes/stacktraces on the web
 - We often index “meta-data” separately
 - Creation date, format, etc.

Tokenization: language issues

- French
 - *L'ensemble* → one token or two?
 - *L ? L' ? Le ?*
 - Want *l'ensemble* to match with *un ensemble*
 - Until at least 2003, it didn't on Google
 - » **Internationalization!**
- German noun compounds are not segmented
 - *Lebensversicherungsgesellschaftsangestellter*
 - 'life insurance company employee'
 - German retrieval systems benefit greatly from a **compound splitter** module
 - Can give a 15% performance boost for German

Tokenization: language issues

- Chinese and Japanese have no spaces between words:
 - 莎拉波娃现在居住在美国东南部的佛罗里达。
 - Not always guaranteed a unique tokenization
- Further complicated in Japanese, with multiple alphabets intermingled
 - Dates/amounts in multiple formats



End-user can express query entirely in Hiragana!

Tokenization: language issues

- Arabic (or Hebrew) is basically written right to left, but with certain items like numbers written left to right
- Words are separated, but letter forms within a word form complex ligatures

← → ← → ← start

استقلت الجزائر في سنة 1962 بعد 132 عام من الاحتلال الفرنسي.

‘Algeria achieved its independence in 1962 after 132 years of French occupation.’

- With Unicode, the order of characters in files matches the conceptual order, and the reversal of displayed characters is handled by the rendering system.

Stop words

- Common words which would appear to be of little value.
 - e.g. the, a, and, to, be
- With a stop list, you exclude from the dictionary entirely the commonest words. Intuition:
 - They have little semantic content
 - There are a lot of them: ~30% of postings for top 30 words
- But the trend is away from doing this:
 - Good compression techniques means the space for including stop words in a system is very small
 - Good query optimization techniques mean you pay little at query time for including stop words.
 - You need them for:
 - Phrase queries: “King of Denmark”
 - Various song titles, etc.: “Let it be”, “To be or not to be”
 - “Relational” queries: “flights to London”

Normalization to terms

- We may need to “normalize” words in indexed text as well as query words into the same form
 - We want to match ***U.S.A.*** and ***USA***
- Result is terms: a **term** is a (normalized) word type, which is an entry in our IR system dictionary
- We most commonly implicitly define equivalence classes of terms by, e.g.,
 - deleting periods to form a term
 - ***U.S.A., USA***
 - deleting hyphens to form a term
 - ***anti-discriminatory, antidiscriminatory***

Normalization: other languages

- Normalization of things like date forms
 - *7月30日 vs. 7/30*
 - *Japanese use of kana vs. Chinese characters*
- Tokenization and normalization may depend on the language and so is intertwined with language detection
- Crucial: Need to “normalize” indexed text as well as query terms **identically**

Case folding

- Reduce all letters to lower case
 - exception: upper case in mid-sentence?
 - e.g., General Motors
 - Fed vs. fed
 - SAIL vs. sail
 - Often best to lower case everything, since users will use lowercase regardless of ‘correct’ capitalization...
- Longstanding Google example:
 - Query: C.A.T.
 - #1 result is for “cats”, not Caterpillar Inc.

Lemmatization

- Reduce inflectional/variant forms to base form
e.g.,
 - *am, are, is* → *be*
 - *car, cars, car's, cars'* → *car*
- *the boy's cars are different colors* → *the boy car be different color*
- Lemmatization implies doing “proper” reduction to dictionary headword form

Stemming

- Reduce terms to their “roots” before indexing
- “Stemming” suggests crude affix chopping
 - language dependent
 - e.g., *automate(s)*, *automatic*, *automation* all reduced to *automat*.

for example compressed and compression are both accepted as equivalent to compress.



for exampl compress and compress ar both accept as equival to compress

Porter's algorithm

- Commonest algorithm for stemming English
 - Results suggest it's at least as good as other stemming options
- Conventions + 5 phases of reductions
 - phases applied sequentially
 - each phase consists of a set of commands
 - sample convention: *Of the rules in a compound command, select the one that applies to the longest suffix.*

Typical rules in Porter

- *sses* → *ss*
- *ies* → *i*
- *ational* → *ate*
- *tional* → *tion*

- Weight of word sensitive rules
- *(m>1) EMENT* →
 - *replacement* → *replac*
 - *cement* → *cement*

Does stemming help?

- English: very mixed results. Helps recall for some queries but harms precision on others
e.g., operative (dentistry) \Rightarrow oper
- Definitely useful for Spanish, German, Finnish,
...
 - 30% performance gains for Finnish!