



# Ανάκληση Πληροφορίας

## Information Retrieval

Διδάσκων –  
Δημήτριος Κατσαρός



# Dictionary και Postings



# Recall basic indexing pipeline

Documents to be indexed.



Friends, Romans, countrymen.  
⋮

Tokenizer

Token stream.

Friends

Romans

Countrymen

Linguistic modules

Modified tokens.

friend

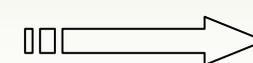
roman

countryman

Indexer

Inverted index.

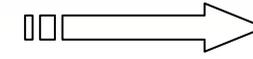
**friend**



2

4

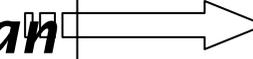
**roman**



1

2

**countryman**



13

16



## Parsing a document

- What format is it in?
  - pdf/word/excel/html?
- What language is it in?
- What character set is in use?

Each of these is a classification problem, which we will study later in the course.

But these tasks are often done heuristically ...



## Complications: Format/language

- Documents being indexed can include docs from many different languages
  - A single index may have to contain terms of several languages.
- Sometimes a document or its components can contain multiple languages/formats
  - French email with a German pdf attachment.
- What is a unit document?
  - A file?
  - An email? (Perhaps one of many in an mbox.)
  - An email with 5 attachments?
  - A group of files (PPT or LaTeX in HTML)



# Tokenization



# Tokenization

- Input: “*Friends, Romans and Countrymen*”
- Output: Tokens
  - *Friends*
  - *Romans*
  - *Countrymen*
- 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? Finlands? Finland's?*
  - *Hewlett-Packard* → *Hewlett* and  
*Packard* as two tokens?
    - *State-of-the-art*: break up hyphenated sequence.
    - *co-education* ?
    - *the hold-him-back-and-drag-him-away-maneuver* ?
    - It's 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

- *3/12/91* *Mar. 12, 1991*
- *55 B.C.*
- *B-52*
- *My PGP key is 324a3df234cb23e*
- *100.2.86.144*
  - Often, don't index as text.
    - But often very useful: think about things like looking up error codes/stacktraces on the web
    - (One answer is using n-grams: Lecture 3)
  - Will often index “meta-data” separately
    - Creation date, format, etc.



# Tokenization: Language issues

- *L'ensemble* → one token or two?
  - *L ? L' ? Le ?*
  - Want *l'ensemble* to match with *un ensemble*
- German noun compounds are not segmented
  - Lebensversicherungsgesellschaftsangestellter
  - 'life insurance company employee'



# 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
- استقلت الجزائر في سنة 1962 بعد 132 عاما من الاحتلال الفرنسي.
- ← → ← → ← start
- ‘Algeria achieved its independence in 1962 after 132 years of French occupation.’
- With Unicode, the surface presentation is complex, but the stored form is straightforward



# Normalization

- Need to “normalize” terms in indexed text as well as query terms into the same form
  - We want to match *U.S.A.* and *USA*
- We most commonly implicitly define equivalence classes of terms
  - e.g., by deleting periods in a term
- Alternative is to do asymmetric expansion:
  - Enter: *window*    Search: *window, windows*
  - Enter: *windows*    Search: *Windows, windows*
  - Enter: *Windows*    Search: *Windows*
- Potentially more powerful, but less efficient



# Normalization: other languages

- Accents: *résumé* vs. *resume*.
- Most important criterion:
  - How are your users like to write their queries for these words?
- Even in languages that standardly have accents, users often may not type them
- German: Tuebingen vs. Tübingen
  - Should be equivalent



## Normalization: other languages

- Need to “normalize” indexed text as well as query terms into the same form

**7 30 vs. 7/30**

- Character-level alphabet detection and conversion
  - Tokenization not separable from this.
  - Sometimes ambiguous:

*Morgen will ich in MIT...*

Is this  
German “mit”?



# 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...



# Stop words

- With a stop list, you exclude from dictionary entirely the commonest words. Intuition:
  - They have little semantic content: *the, a, and, to, be*
  - They take a lot of space: ~30% of postings for top 30
- But the trend is away from doing this:
  - Good compression techniques (lecture 5) means the space for including stopwords 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”



# Thesauri and soundex

- Handle synonyms and homonyms
  - Hand-constructed equivalence classes
    - e.g., *car* = *automobile*
    - *color* = *colour*
- Rewrite to form equivalence classes
- Index such equivalences
  - When the document contains *automobile*, index it under *car* as well (usually, also vice-versa)
- Or expand query?
  - When the query contains *automobile*, look under *car* as well



# Soundex

- Traditional class of heuristics to expand a query into phonetic equivalents
  - Language specific – mainly for names
  - E.g., *chebyshev* → *tchebycheff*
- More on this later ...



# 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” suggest 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 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*



## Other stemmers

- Other stemmers exist, e.g., Lovins stemmer  
<http://www.comp.lancs.ac.uk/computing/research/stemming/general/lovins.htm>
  - Single-pass, longest suffix removal (about 250 rules)
  - Motivated by linguistics as well as IR
- Full morphological analysis – at most modest benefits for retrieval
- Do stemming and other normalizations help?
  - Often very mixed results: really help recall for some queries but harm precision on others



# Language-specificity

- Many of the above features embody transformations that are
  - Language-specific and
  - Often, application-specific
- These are “plug-in” addenda to the indexing process
- Both open source and commercial plug-ins available for handling these



## Dictionary entries – first cut

<i>ensemble.french</i>
<i>.japanese</i>
<i>MIT.english</i>
<i>mit.german</i>
<i>guaranteed.english</i>
<i>entries.english</i>
<i>sometimes.english</i>
<i>tokenization.english</i>

These may be grouped by language (or not...).

More on this in ranking/query processing.