

# Introduction to **Information Retrieval**

Term vocabulary and postings lists –  
preprocessing steps

# Documents

---

- Last lecture: Simple Boolean retrieval system
- Our assumptions were:
  - We know what a document is.
  - We can “machine-read” each document.
- This can be complex in reality.

# Parsing a document

---

- Convert byte sequence into a linear sequence of characters
- Requirements
  - Deal with format and language of each document
  - What is the encoding? E.g., UTF-8
  - What format is it in? pdf, word, excel, html, etc.
  - What language is it in?
  - What character set is in use?
- Each of these is a classification problem
- Alternative: use heuristics

# Format/Language: Complications

---

- A single index usually contains terms of several languages.
- A document may contain multiple languages/formats
  - French email with Spanish pdf attachment
  - Code switching in social media

# Format/Language: Complications

---

- What is the document unit for indexing?
  - A file?
  - An email? An email file can contain a sequence of messages; each message can be considered a document
  - An email with 5 attachments?
  - Multiple files may be combined into one document (ppt or latex in HTML)
- Upshot: Answering the question “what is a document?” is not trivial and requires some design decisions.

# Definitions

---

- **Word** – A delimited string of characters as it appears in the text
- **Term** – A “normalized” word (case, morphology, spelling etc); actually an equivalence class of words; usually what is included in an IR system’s dictionary
- **Token** – An instance of a word or term occurring in a document.
- **Type** – The same as a term in most cases: an equivalence class of tokens.

# Recall: Inverted index construction

---

- Input:

Friends, Romans, countrymen. So let it be with Caesar ...

- Output:

friend roman countryman so ...

- Each token is a candidate for a postings entry.
- What are valid tokens to emit?

# Exercises

---

*In June, the dog likes to chase the cat in the barn.*

– How many word tokens? How many word types?

Why tokenization is difficult even in English?

**Tokenize:** *Mr. O'Neill thinks that the boys' stories about Chile's capital aren't amusing.*



# Tokenization problems: One word or two? (or several)

---

- Hewlett-Packard
- State-of-the-art
- co-education
- the hold-him-back-and-drag-him-away maneuver
- data base
- San Francisco
- Los Angeles-based company
- cheap San Francisco-Los Angeles fares
- York University vs. New York University

# Tokenization problems: Numbers

---

- 3/20/91
- 20/3/91
- Mar 20, 1991
- B-52
- 100.2.86.144
- (800) 234-2333
- 800.234.2333
- Older IR systems may not index numbers . . .
- . . . but generally it's a useful feature.

## Problems in tokenization for other languages, e.g., no whitespace in Chinese

---

莎拉波娃现在居住在美国东南部的佛罗里达。今年4月9日，莎拉波娃在美国第一大城市纽约度过了18岁生日。生日派对上，莎拉波娃露出了甜美的微笑。

## Ambiguous segmentation in Chinese

---

和尚

The two characters can be treated as one word meaning 'monk' or as a sequence of two words meaning 'and' and 'still'.

## Other cases of “no whitespace”

---

- Compounds in Dutch, German, Swedish
- Computerlinguistik → Computer + Linguistik
- Lebensversicherungsgesellschaftsangestellter
- → leben + versicherung + gesellschaft + angestellter
- Inuit: tusaatsiarunnanngittualuujunga (I can't hear very well.)
- Many other languages with segmentation difficulties: Finnish, Urdu, . . .
- **CamelCase in social media**

# Japanese

---

ノーベル平和賞を受賞したワンガリ・マータイさんが名誉会長を務めるMOTTAINAIキャンペーンの一環として、毎日新聞社とマガジンハウスは「私の、もったいない」を募集します。皆様が日ごろ「もったいない」と感じて実践していることや、それに関わるエピソードを800字以内の文章にまとめ、簡単な写真、イラスト、図などを添えて10月20日までにお送りください。大賞受賞者には、50万円相当の旅行券とエコ製品2点の副賞が贈られます。

4 different “alphabets”: Chinese characters, hiragana syllabary for inflectional endings and functional words, katakana syllabary for transcription of foreign words and other uses, and latin. No spaces (as in Chinese). End user can express query entirely in hiragana!

## Arabic script

---

كِتَابٌ ← كِتَابٌ  
un b ā t i k  
/kitābun/ 'a book'

## Arabic script: Bidirectionality

---

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

← →      ← →      ← START

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

Bidirectionality is not a problem if text is coded in Unicode.



# Accents and diacritics

---

- Accents: résumé vs. resume (simple omission of accent)
- Umlauts: Universität vs. Universitaet (substitution with special letter sequence “ae”)
- Most important criterion: How are users likely to write their queries for these words?
- Even in languages that standardly have accents, users often do not type them. (Polish?)

# Case folding

---

- Usually: Reduce all letters to lower case
- Possible exceptions: capitalized words in mid-sentence
- MIT vs. mit
- Fed vs. fed
- It's often best to lowercase everything since users will use lowercase regardless of correct capitalization

# Normalization

---

- Need to “normalize” terms in indexed text as well as query terms into the same form.
- Example: We want to match *U.S.A.* and *USA*
- We commonly implicitly define **equivalence classes of terms**
  - Can use hand-constructed rules, e.g., ‘car’ & ‘automobile’
- Alternatively: do asymmetric expansion
  - window → window, windows
  - windows → Windows, windows
  - Windows (no expansion)
- More powerful, but less efficient
- Why don’t you want to put *window*, *Window*, *windows*, and *Windows* in the same equivalence class?

# Normalization: complex for multiple languages

---

- Normalization and language detection interact.
- *PETER WILL NICHT MIT.* → MIT = mit
- *He got his PhD from MIT.* → MIT ≠ mit

# Stop words

---

- stop words: **extremely common words** which would appear to be of little value in helping select documents matching a user need
- Examples: *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*
- Stop word elimination used to be standard in older IR systems.
- But you need stop words for phrase queries, e.g. “King of Denmark”
- **Most web search engines index stop words.**

# Stemming and Lemmatization

---

- Goal of both same: reduce inflectional forms and derivationally related forms to a common base form
- Lemmatization implies doing “proper” reduction to dictionary headword form (the **lemma**), using dictionary and morphological analysis of words
- Stemming refers to a heuristic process that chops off the ends of words in the hope of achieving the goal correctly most of the time

# Lemmatization

---

- Reduce inflectional/variant forms to base form
- Example: *am, are, is* → *be*
- Example: *car, cars, car's, cars'* → *car*
- Example: *the boy's cars are different colors* → *the boy car be different color*
- Inflectional morphology (*cutting* → *cut*) vs. derivational morphology (*destruction* → *destroy*)

# Stemming

---

- Heuristic process that **chops off the ends of words** in the hope of achieving what “principled” lemmatization attempts to do with a lot of linguistic knowledge.
- Language dependent
- Often inflectional **and** derivational
- Example for derivational: *automate, automatic, automation* all reduce to *automat*



# Porter algorithm

---

- Most common algorithm for stemming English
- Results suggest that it is at least as good as other stemming options
- Conventions + 5 phases of reductions
- Phases are applied sequentially
- Each phase consists of a set of commands.
  - Sample command: Delete final *ement* if what remains is longer than 1 character
  - replacement → replac
  - cement → cement

## Porter stemmer: A few rules

---

### Rule

SSES → SS

IES → I

SS → SS

S →

### Example

caresses → caress

ponies → poni

caress → caress

cats → cat

Sample convention: Of the rules in a compound command, select the one that applies to the longest suffix.

# Three stemmers: A comparison

---

*Sample text:* Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation

*Porter stemmer:* such an analysi can reveal featur that ar not easili visibl from the variat in the individu gene and can lead to pictur of express that is more biolog transpar and access to interpret

*Lovins stemmer:* such an analys can reve featur that ar not eas vis from th vari in th individu gen and can lead to a pictur of expres that is mor biolog transpar and acces to interpres

*Paice stemmer:* such an analys can rev feat that are not easy vis from the vary in the individ gen and can lead to a pict of express that is mor biolog transp and access to interpret

# Does stemming improve effectiveness?

---

- In general, stemming increases effectiveness for some queries, and decreases effectiveness for others.
- Queries where stemming is likely to help: [tartan sweaters], [sightseeing tour san francisco] (equivalence classes: {sweater, sweaters}, {tour, tours})
- Porter Stemmer equivalence class *oper* contains all of *operate operating operates operation operative operatives operational*.
- Queries where stemming hurts: [operational AND research], [operating AND system], [operative AND dentistry]