

## Information Retrieval

INFO 4300 / CS 4300

---

- Last classes
  - Text acquisition
    - » Web crawlers

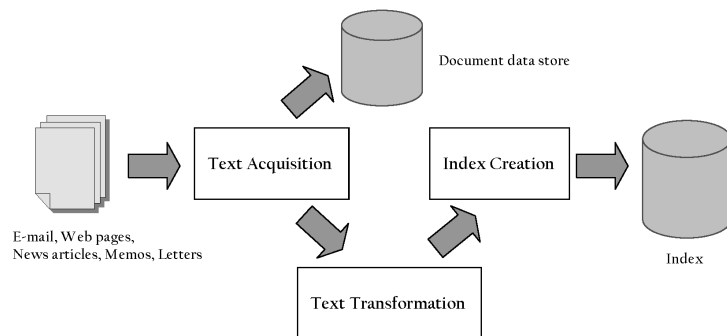
## Topics for Today

---

- Text transformation
  - Word occurrence statistics
  - Tokenizing
  - Stopping and stemming

## Indexing Process

---



## Pop Quiz: Indexing

---

Doc 1: King Arthur : I am your king  
Doc 2: Woman : Well I didn't vote for you  
Doc 3: King Arthur : You don't vote for kings  
Doc 4: Woman : Well how'd you become king then

Show the inverted index for any word in the document.

## Processing Text

---

- Converting documents to a more consistent set of *index terms*
- Why?
  - Sometimes not clear where words begin and end
    - » Not even clear what a word is in some languages
      - ◆ e.g., Chinese, Korean
  - Matching the exact string of characters typed by the user is too restrictive
    - » i.e., it doesn't work very well in terms of effectiveness
  - Not all words are of equal value in a search

## Topics for Today

---

- Text transformation
  - ➔ – Word occurrence statistics
    - » Distribution
    - » Vocabulary growth
  - Tokenizing
  - Stopping and stemming

## Text Statistics

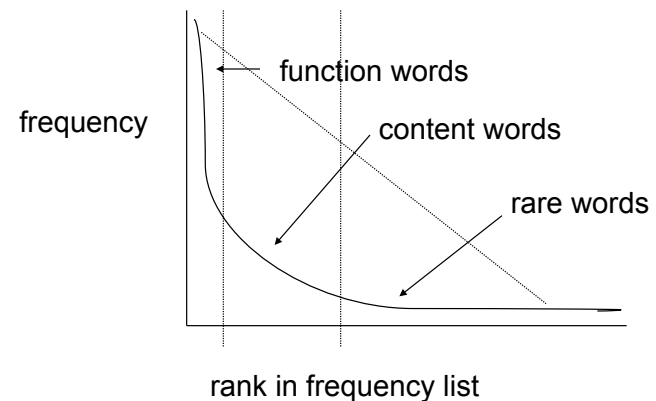
---

- Many statistical characteristics of word occurrences are predictable
- Retrieval models and ranking algorithms depend heavily on them
  - e.g., important words occur often in documents but are not high frequency in collection [Luhn, 1958]

## Distribution of word frequencies

---

- is very *skewed*



## Zipf's Law

- Distribution of word frequencies is very *skewed*
  - a few words occur very often, many words hardly ever occur
  - e.g., two most common words (“the”, “of”) make up about 10% of all word occurrences in text documents
- Zipf's “law”:
  - observation that rank ( $r$ ) of a word times its frequency ( $f$ ) is approximately a constant ( $k$ )
    - » assuming words are ranked in order of decreasing frequency
  - i.e.,  $r \cdot f \approx k$  or  $r \cdot P_r \approx c$ , where  $P_r$  is probability of word occurrence for the  $r$ th ranked word and  $c \approx 0.1$  for English

## Zipf's Law (*Tom Sawyer*)

Word	Freq. ( $f$ )	Rank ( $r$ )	$f \cdot r$	Word	Freq. ( $f$ )	Rank ( $r$ )	$f \cdot r$
the	3332	1	3332	turned	51	200	10200
and	2972	2	5944	you'll	30	300	9000
a	1775	3	5235	name	21	400	8400
he	877	10	8770	comes	16	500	8000
but	410	20	8400	group	13	600	7800
be	294	30	8820	lead	11	700	7700
there	222	40	8880	friends	10	800	8000
one	172	50	8600	begin	9	900	8100
about	158	60	9480	family	8	1000	8000
more	138	70	9660	brushed	4	2000	8000
never	124	80	9920	sins	2	3000	6000
Oh	116	90	10440	Could	2	4000	8000
two	104	100	10400	Applausive	1	8000	8000

Manning and Schutze SNLP

## Zipf's Law

- Useful as a rough description of the frequency distribution of words in human languages
- Behavior occurs in a surprising variety of situations
  - References to scientific papers
  - Web page in-degrees, out-degrees
  - Royalties to pop-music composers

## News Collection (AP89) Statistics

Total documents	84,678
Total word occurrences	39,749,179
Vocabulary size	198,763
Words occurring > 1000 times	4,169
Words occurring once	70,064

## Top 50 Words from AP89

Word	Freq.	r	Pr(%)	r.Pr	Word	Freq.	r	Pr(%)	r.Pr
the	2,420,778	1	6.49	0.065	has	136,007	26	0.37	0.095
of	1,045,733	2	2.80	0.056	are	130,322	27	0.35	0.094
to	968,882	3	2.60	0.078	not	127,493	28	0.34	0.096
a	892,429	4	2.39	0.096	who	116,364	29	0.31	0.090
and	865,644	5	2.32	0.120	they	111,024	30	0.30	0.089
in	847,825	6	2.27	0.140	its	111,021	31	0.30	0.092
said	504,593	7	1.35	0.095	had	103,943	32	0.28	0.089
for	363,865	8	0.98	0.078	will	102,949	33	0.28	0.091
that	347,072	9	0.93	0.084	would	99,503	34	0.27	0.091
was	293,027	10	0.79	0.079	about	92,983	35	0.25	0.087
on	291,947	11	0.78	0.086	i	92,005	36	0.25	0.089
he	250,919	12	0.67	0.081	been	88,786	37	0.24	0.088
is	245,843	13	0.65	0.086	this	87,286	38	0.23	0.089
with	223,846	14	0.60	0.084	their	84,638	39	0.23	0.089
at	210,064	15	0.56	0.085	new	83,449	40	0.22	0.090
by	209,586	16	0.56	0.090	or	81,796	41	0.22	0.090
it	195,621	17	0.52	0.089	which	80,385	42	0.22	0.091
from	189,451	18	0.51	0.091	we	80,245	43	0.22	0.093
as	181,714	19	0.49	0.093	more	76,388	44	0.21	0.090
be	157,300	20	0.42	0.084	after	75,165	45	0.20	0.091
were	153,913	21	0.41	0.087	us	72,045	46	0.19	0.089
an	152,576	22	0.41	0.090	percent	71,956	47	0.19	0.091
have	149,749	23	0.40	0.092	up	71,082	48	0.19	0.092
his	142,285	24	0.38	0.092	one	70,266	49	0.19	0.092
but	140,880	25	0.38	0.094	people	68,988	50	0.19	0.093

## Lower frequency words from AP89

Word	Freq.	r	Pr(%)	r*Pr
assistant	5,095	1,021	.013	0.13
sewers	100	17,110	$2.56 \times 10^{-4}$	0.04
toothbrush	10	51,555	$2.56 \times 10^{-5}$	0.01
hazmat	1	166,945	$2.56 \times 10^{-6}$	0.04

## Topics for Today

- Text transformation
  - Word occurrence statistics
    - » Distribution
    - » Vocabulary growth
  - Tokenizing
  - Stopping and stemming

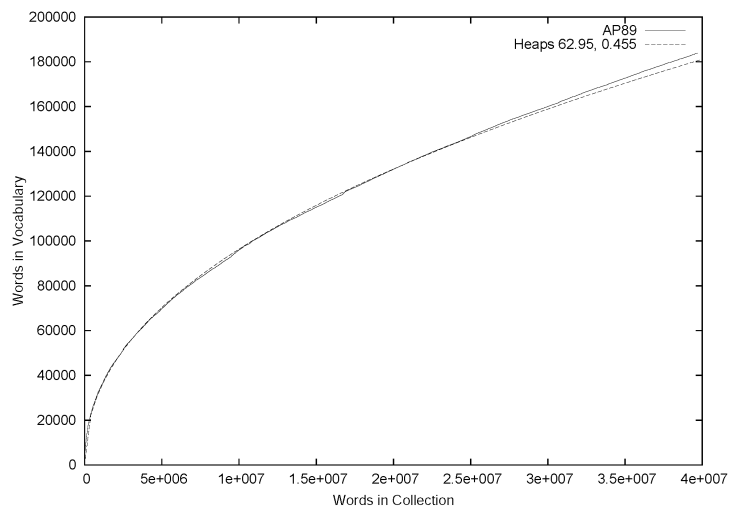
## Vocabulary Growth

- As corpus grows, so does vocabulary size
  - Fewer new words when corpus is already large
- Observed relationship (*Heaps' Law*):

$$v = k * n^\beta$$

$v$  is vocabulary size (number of unique words),  
 $n$  is the number of words in corpus,  
 $k, \beta$  are parameters that vary for each corpus  
 (typical values given are  $10 \leq k \leq 100$  and  $\beta \approx 0.5$ )

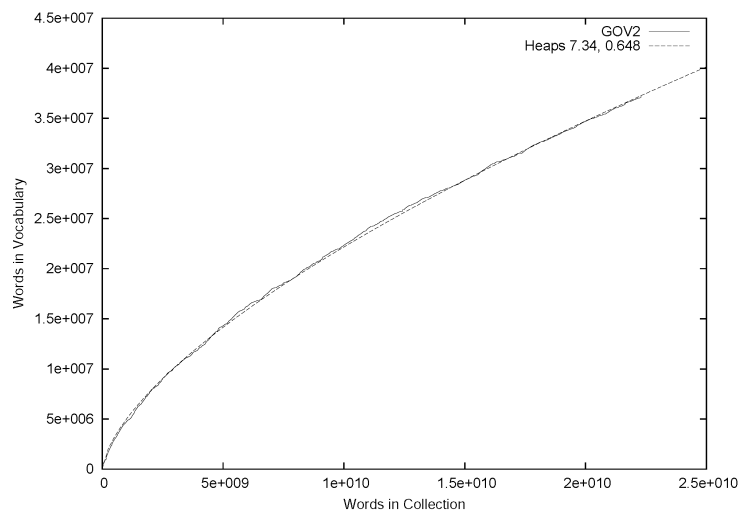
## Heaps' Law: AP89 Example



## Heaps' Law Predictions

- Predictions for TREC collections are accurate for large numbers of words
  - e.g., first 10,879,522 words of the AP89 collection scanned
  - prediction is 100,151 unique words
  - actual number is 100,024
- Predictions for small corpora (i.e. < 1000 words) are much worse

## GOV2 (Web) Example



## Web Example

- Heaps' Law works with very large corpora
  - new words occurring even after seeing 30 million!
  - parameter values different than typical TREC values
- New words come from a variety of sources
  - » spelling errors, invented words (e.g. product, company names), code, other languages, email addresses, etc.
- Search engines must deal with these large and growing vocabularies

## Topics for Today

---

- Text transformation
  - Word occurrence statistics
- ➔ – Tokenizing
- Stopping and stemming

## Tokenizing

---

- Forming words from sequence of characters
- Surprisingly complex in English, can be harder in other languages
- Early IR systems:
  - any sequence of alphanumeric characters of length 3 or more
  - terminated by a space or other special character
  - upper-case changed to lower-case

## Tokenizing

---

- Example:
  - “Bigcorp's 2007 bi-annual report showed profits rose 10%.” becomes
  - “bigcorp 2007 annual report showed profits rose”
- Too simple for most search applications  
Why? Too much information lost
  - Small decisions in tokenizing can have major impact on effectiveness of some queries

## Tokenizing Problems

---

- Small words can be important in some queries, usually in combinations
  - » xp, ma, pm, ben e king, el paso, master p, gm, j lo, world war II
- Both hyphenated and non-hyphenated forms of many words are common
  - Sometimes hyphen is not needed
    - » e-bay, wal-mart, active-x, cd-rom, t-shirts
  - At other times, hyphens should be considered either as part of the word or a word separator
    - » winston-salem, mazda rx-7, e-cards, pre-diabetes, t-mobile, spanish-speaking

## Tokenizing Problems

---

- Special characters are an important part of tags, URLs, code in documents
- Capitalized words can have different meaning from lower case words
  - Bush, Apple
- Apostrophes can be a part of a word, a part of a possessive, or just a mistake
  - rosie o'donnell, can't, don't, 80's, 1890's, men's straw hats, master's degree, england's ten largest cities, shriner's

## Tokenizing Problems

---

- Numbers can be important, including decimals
  - nokia 3250, top 10 courses, united 93, quicktime 6.5 pro, 92.3 the beat
- Periods can occur in numbers, abbreviations, URLs, ends of sentences, and other situations
  - I.B.M., Ph.D., cs.umass.edu, F.E.A.R.
- Note: tokenizing steps for queries must be identical to steps for documents

## Tokenizing Process

---

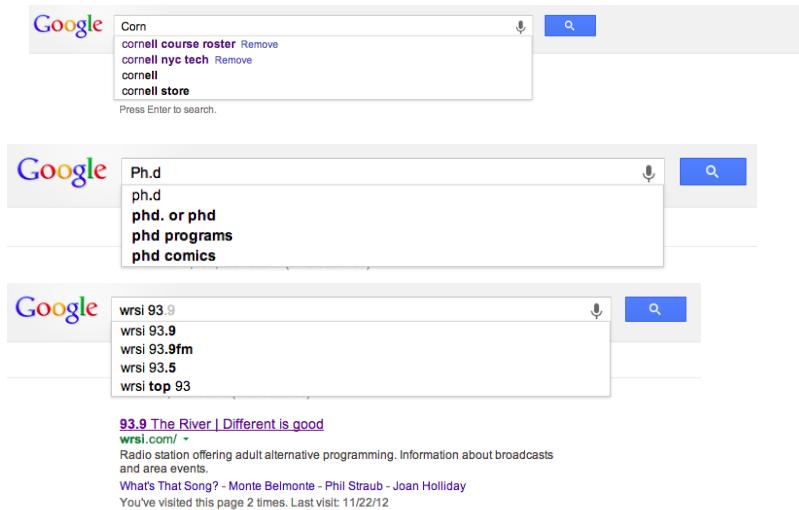
- First step is to use parser to identify appropriate parts of document to tokenize
- Defer complex decisions to other components
  - word is any sequence of alphanumeric characters, terminated by a space or special character, with everything converted to lower-case
  - everything indexed
  - example: 92.3 → 92 3 but search finds documents with 92 and 3 adjacent
  - incorporate some rules to reduce dependence on query transformation components

## Tokenizing Process

---

- Not that different than simple tokenizing process used in past
- Examples of rules used with TREC
  - Apostrophes in words ignored
    - » o'connor → oconnor bob's → bobs
  - Periods in abbreviations ignored
    - » I.B.M. → ibm Ph.D. → phd

## What does Google do?



## Topics for Today

- Text transformation
  - Word occurrence statistics
  - Tokenizing
  - ➔ – Stopping and stemming

## Stopping

- Function words (determiners, prepositions) have little meaning on their own
- High occurrence frequencies
- Treated as *stopwords* (i.e. removed)
  - reduce index space, improve response time, improve effectiveness
- Can be important in combinations
  - e.g., “to be or not to be”

## Stopping

- Stopword list can be created from high-frequency words or based on a standard list
- Lists are customized for applications, domains, and even parts of documents
  - e.g., “click” is a good stopwords for anchor text
- **Best policy** is to index all words in documents, make decisions about which words to use at query time