

Introduction to **Information Retrieval**

Lectures 7: Dictionaries and tolerant retrieval

Dictionary

- Given a query and an inverted index, first tasks:
 - Determine whether each query term exists in the vocabulary
 - If yes, identify the pointers to the corresponding postings lists
- The vocabulary lookup operation uses a data structure called dictionary.

A naïve dictionary

- An array of struct:

term	document frequency	pointer to postings list
a	656,265	→
aachen	65	→
...
zulu	221	→

char[20] int Postings *

20 bytes 4/8 bytes 4/8 bytes

- How do we store a dictionary in memory efficiently?
- How do we quickly look up elements at query time?

Dictionary data structures

- Two main choices:
 - Hashtables
 - Trees
- Some IR systems use hashtables, some trees

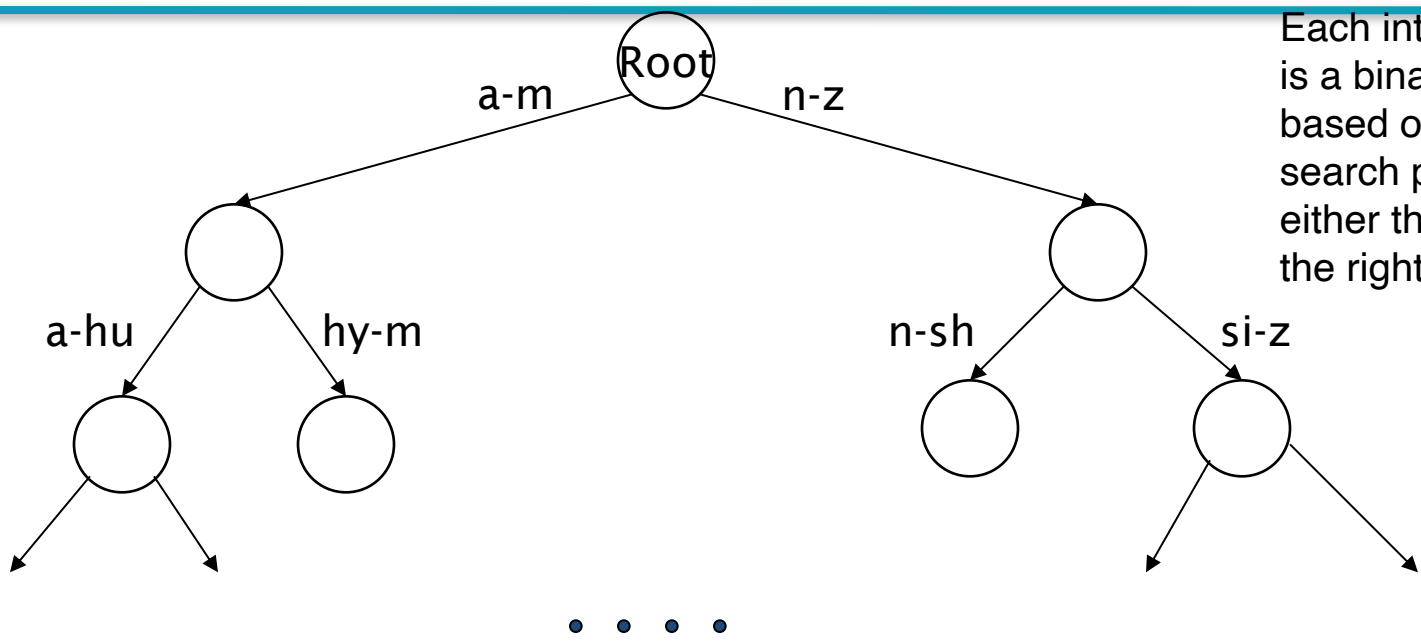
Hashtables

- Each vocabulary term is hashed to an integer over a large enough space where collisions are unlikely
 - (We assume you've seen hashtables before)
- Pros:
 - Lookup is faster than for a tree: $O(1)$
- Cons:
 - No easy way to find minor variants (judgment/judgement) since they could be hashed to very different indexes:
 - No easy way to handle wild-card queries / prefix search
 - If vocabulary keeps growing, need to occasionally do the expensive operation of rehashing *everything*

Trees

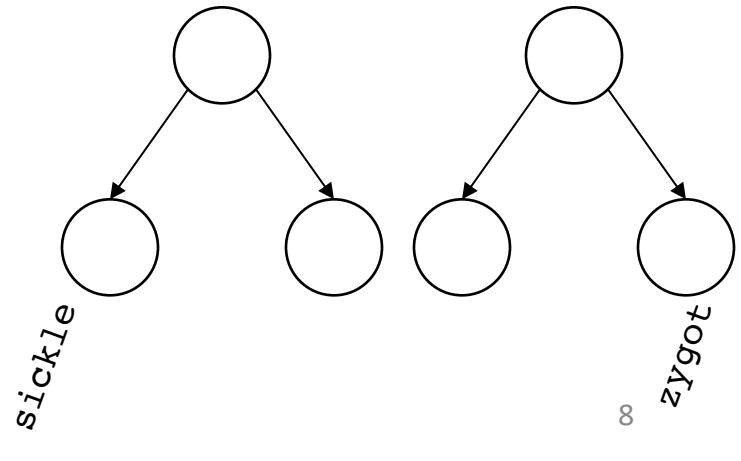
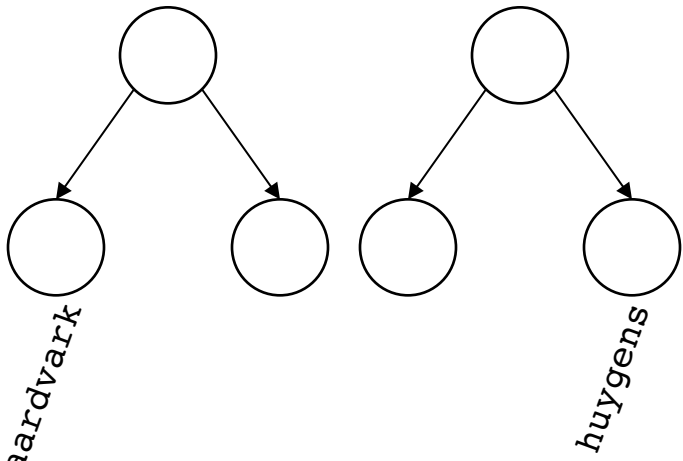
- Simplest: binary tree
- More usual: B-trees

Tree: binary tree

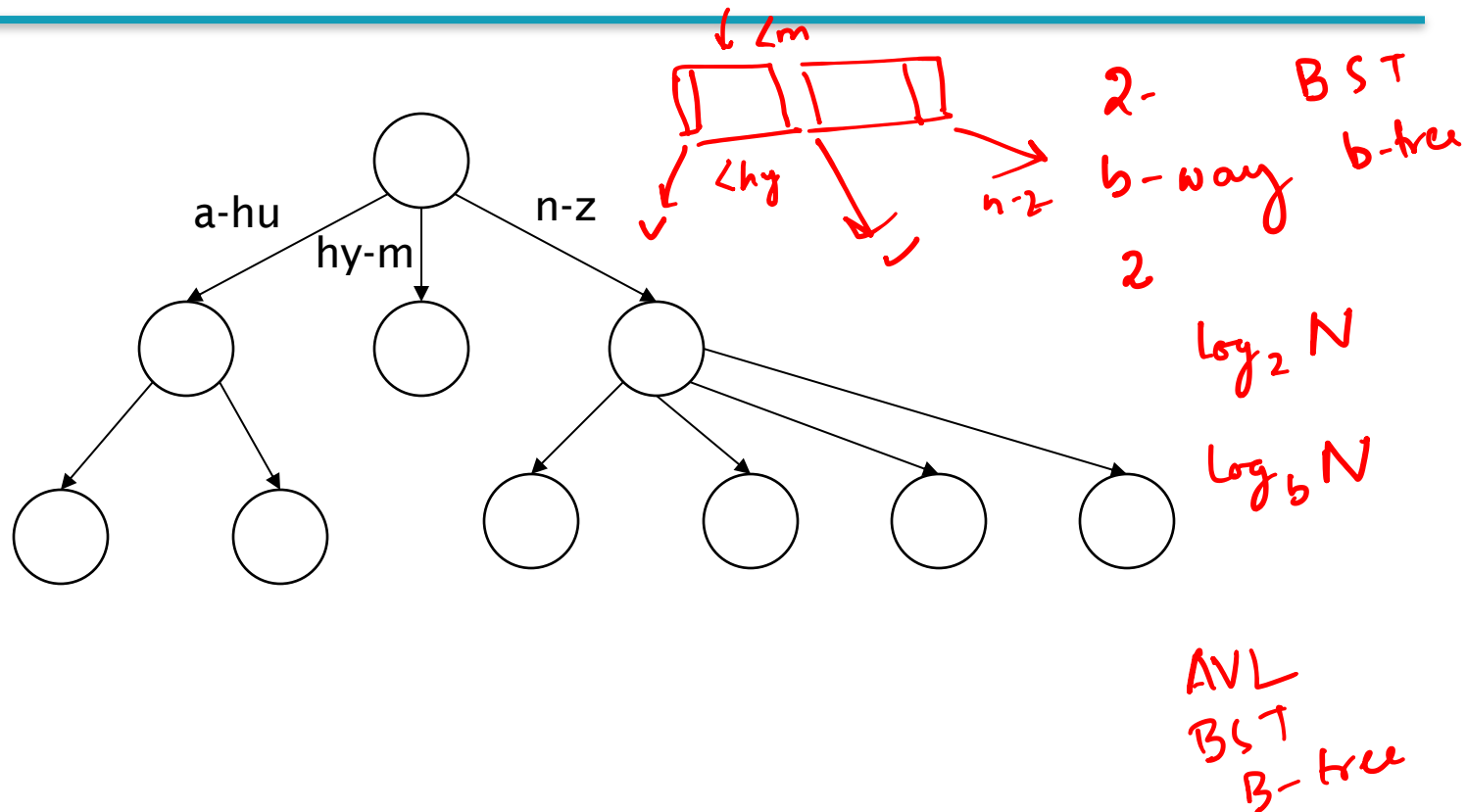


Each internal node is a binary test, based on which the search proceeds to either the left or the right subtree

...



Tree: B-tree



- Definition: Every internal node has a number of children in the interval $[a, b]$ where a, b are appropriate natural numbers, e.g., $[2, 4]$.

Trees

- Trees require a standard ordering of characters and hence strings ... but we typically have one
- Pros:
 - Solves the prefix problem (terms starting with *hyp*)
- Cons:
 - Slower: $O(\log M)$ [and this requires *balanced* tree]
 - Rebalancing binary trees is expensive
 - But B-trees mitigate the rebalancing problem

WILD-CARD QUERIES

Dictionary data structures Summary

- Two main choices: Hashtables, Trees
- Some IR systems use hashtables, some trees
- Hashtables:
 - Pros: $O(1)$ for large enough table space.
 - Cons: 1) minor mistakes in query, 2) wildcard search hard. 3) as vocab grows, you may need to rehash everything.
- Trees:
 - Simplest: binary Search Tree
 - More general/useful: B-Tree, solves prefix problem. $O(\log M)$

Wild-card queries: *

- ***mon****: find docs having any word beginning with “mon”
- Why wild-card queries?
 - Users may not be sure of a spelling
 - Users may want to match multiple variants of a term, e.g., query judicia* for matching both judicial and judiciary

mon
jury*

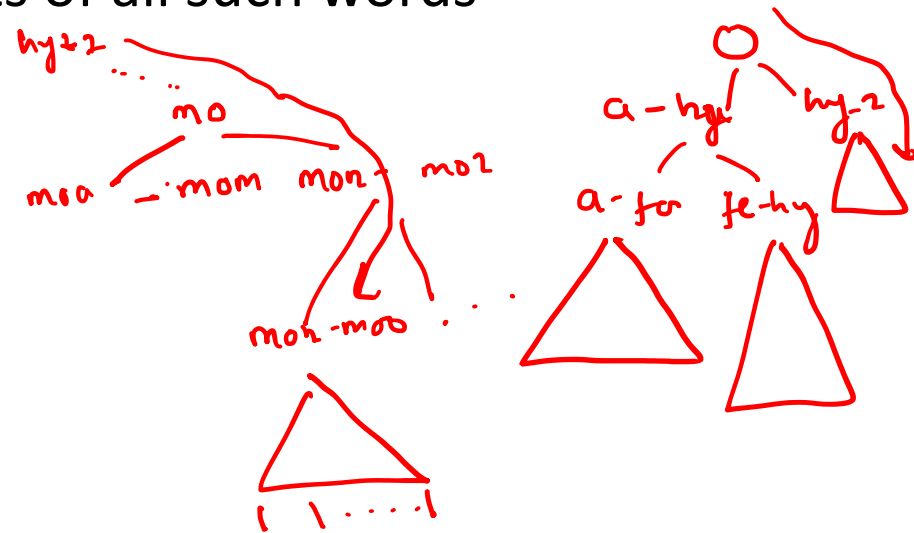
- Types of wild-card queries

mon challenges*

- Trailing wild-card queries, e.g., *mon**
 - Leading wild-card queries, e.g., **mon*
 - General wild-card queries, e.g., *s*dney, fas*in*te*
- **Key challenge: which dictionary terms match a wild-card query?**

Trailing Wild-card queries

- Trailing wild-card queries (mon^*) easier
 - Use a binary tree (or B-tree) over the dictionary:
 - Retrieve all dictionary words in range: $\text{mon} \leq w < \text{moo}$
 - Then process postings lists of all such words



Leading Wild-card queries

- ***mon**: find words ending in “mon”: slightly harder than trailing wild-card queries
- Maintain an additional **B-tree for terms backwards (reverse B-tree)**
 - Each root-to-leaf path corresponds to a term in the dictionary written backwards
 - Can retrieve all words in range: **nom** ≤ **w** < **non**.

Query processing

- At this point, we have an enumeration of all terms in the dictionary that match the wild-card query.
- We still have to look up the postings for each enumerated term.

Exercise: from what we have discussed, how can we enumerate all terms matching the query **pro*cent** ?

B-trees handle *'s at the end of a query term


- How can we handle *'s in the middle of query term?
 - *co*tion*
- We could look up *co** AND **tion* in a B-tree and intersect the two term sets
 - Expensive
- The solution: transform wild-card queries so that the *'s occur at the end
 - This gives rise to the **Permuterm Index**.

Permuterm index

- Introduce \$ as a special symbol to mark the end of a term (a symbol that does not appear in the text)
- For term **hello**, index under:
 - **hello\$, ello\$h, llo\$he, lo\$hel, o\$hell, \$hello**
where \$ is a special symbol.
 - Various rotations of each term (augmented with \$) all link to the original vocabulary term
- Permuterm vocabulary: set of all rotated terms in the permuterm index

Handling wild-card queries with Permuterm index

- Given a wild-card query
 - Rotate the query so that the * symbol appears at the end of the string
 - Look up the rotated query in the permuterm index
- Queries:
 - **X** lookup on **X\$** **X*** lookup on **\$X***
 - ***X** lookup on **X\$*** ***X*** lookup on **X***
 - **X*Y** lookup on **Y\$X*** **X*Y*Z** ??? Exercise!



Query = *hel*o*
X=hel, Y=o
 Lookup *o\$hel**

Queries having multiple wild-cards

- `fi*mo*er`
- First enumerate all dictionary terms that are in the permuterm index of `er$fi*`
- Not all such terms will have “mo” in the middle – need to filter out mismatched terms exhaustively

Permuterm query processing

- Rotate query wild-card to the right
- Now use B-tree lookup as before.
- Once the permuterm index enables us to identify the original vocabulary terms matching a wild-card query, we can look up these terms in the usual way
- *Permuterm problem: \approx quadruples lexicon size*



Empirical observation for English.

k -gram (e.g. Bigram) indexes

- Enumerate all k -grams (sequence of k chars) occurring in any term
- Use \$ as a special character to denote the beginning and end of each term
- e.g., from text “***April is the cruelest month***” we get the 2-grams (*bigrams*)

\$a,ap,pr,ri,il,l\$, \$i,is,s\$, \$t,th,he,e\$, \$c,cr,ru,
ue,el,le,es,st,t\$, \$m,mo,on,nt,h\$

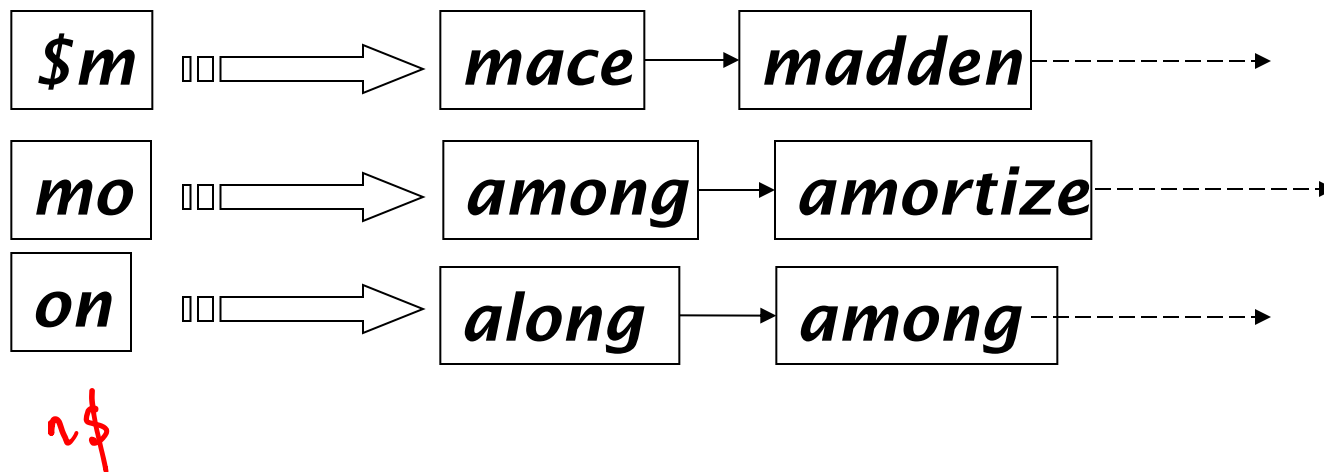
k -gram (e.g. Bigram) indexes

- Dictionary contains all k -grams that occur in any term in the vocabulary
- Maintain a second inverted index from k -grams to dictionary terms that match each k -gram
 - Each postings list points from a k -gram to all vocabulary terms containing that k -gram

april → 1, 2
january → 3, 4

Bigram index example

- The k -gram index finds *terms* based on a query consisting of k -grams (here $k=2$).



Processing wild-cards

- Query *mon** can now be run as
 - *\$m AND mo AND on*
- Gets terms that match AND version of our wildcard query.
- But we'd enumerate *moon*.
- Must **post-filter** these terms against query.
- Surviving enumerated terms are then looked up in the term-document inverted index.
- Fast, space efficient (compared to permuterm).

Processing wild-card queries

- As before, we must execute a Boolean query for each enumerated, filtered term.
- Wild-cards can result in expensive query execution (very large disjunctions...)
 - `pyth*` AND `prog*`
- If you encourage “laziness” people will respond!

Search

Type your search terms, use '*' if you need to.
E.g., `Alex*` will match Alexander.

Recap

- Dictionary Data Structures:
 - Hashtables
 - B-trees

- Tolerant Retrieval
 - Wildcard queries
 - permuterm index
 - k-gram index
 - Spell correction contd ...

SPELLING CORRECTION

Spell correction

- Two principal uses
 - Correcting document(s) being indexed
 - Correcting user queries to retrieve “right” answers
- Two main flavors:
 - Isolated word
 - Check each word on its own for misspelling: ***jacson***
 - Will not catch typos resulting in correctly spelled words
 - e.g., ***from*** → ***form***
 - Context-sensitive
 - Look at surrounding words,
 - e.g., ***I flew form Heathrow to Narita.***

Document correction

- Needed for OCR'ed documents
 - Correction algorithms are tuned for this: rn/m
 - Can use domain-specific knowledge
 - E.g., OCR can confuse O and D more often than it would confuse O and I
 - O and I adjacent on the QWERTY keyboard, so more likely interchanged in typing.
- But often we don't change the documents and instead fix the query-document mapping

Query mis-spellings

- Our principal focus here
 - E.g., the query “*IIT Khargapur*”
- We can either
 - Retrieve documents indexed by the correct spelling, OR
 - Return several suggested alternative queries with the correct spelling
 - *Did you mean ... ?*

Isolated word correction

- Fundamental premise – there is a lexicon from which the correct spellings come
- Two basic choices for this
 - A standard lexicon such as
 - Webster’s English Dictionary
 - An “industry-specific” lexicon – hand-maintained
 - The lexicon of the indexed corpus
 - E.g., all words on the web
 - All names, acronyms etc.
 - (Including the mis-spellings)

Isolated word correction

- Given a lexicon and a character sequence Q , return the words in the lexicon closest to Q
- What's "closest"?
- We'll study several alternatives
 - Edit distance (Levenshtein distance)
 - Weighted edit distance
 - n -gram overlap

Edit distance

- Given two strings S_1 and S_2 , the minimum number of operations to convert one to the other
- Operations are typically character-level
 - Insert, Delete, Replace, (Transposition)
- E.g., the edit distance from **dof** to **dog** is 1
 - From **cat** to **act** is 2 (Just 1 with transpose.)
 - from **cat** to **dog** is 3.
- Generally found by dynamic programming.

cat
 RL ↓
 ca t
 R2 ↓
 a c t

cat
 (d a t
 ↓
 d o t
 ↓
 d o g

C A T
A C T

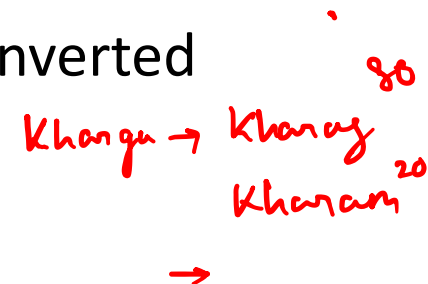
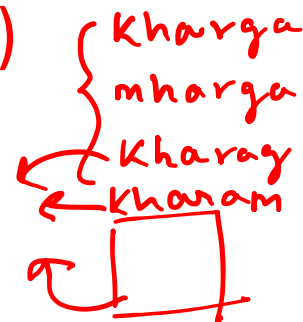
Weighted edit distance

- As above, but the weight of an operation depends on the character(s) involved
 - Meant to capture OCR or keyboard errors
Example: **m** more likely to be mis-typed as **n** than as **q**
 - Therefore, replacing **m** by **n** is a smaller edit distance than by **q**
 - This may be formulated as a probability model
- Requires weight matrix as input
- Modify dynamic programming to handle weights

$26 \times 26 \times 3$
 $m \quad n \quad r -$
 $.$

Using edit distances for correction

- Given query, first enumerate all character sequences within a preset (weighted) edit distance (e.g., 2)
- Intersect this set with list of “correct” words
- Show terms you found to user as suggestions
- Alternatively,
 - We can look up all possible corrections in our inverted index and return all docs ... slow
 - We can run with a single most likely correction
- The alternatives disempower the user, but save a round of interaction with the user



n -gram overlap

- Enumerate all the n -grams in the query string as well as in the lexicon
- Use the n -gram index (recall wild-card search) to **retrieve all lexicon terms matching any of the query n -grams**
- Variations
 - Can threshold by number of matching n -grams
 - Variants – weight by keyboard layout, etc.

Example with trigrams

- Suppose the text is ***november***
 - Trigrams are *nov, ove, vem, emb, mbe, ber.*
- The query is ***december***
 - Trigrams are *dec, ece, cem, emb, mbe, ber.*
- So 3 trigrams overlap (of 6 in each term)
- How can we turn this into a normalized measure of overlap?

One option – Jaccard coefficient

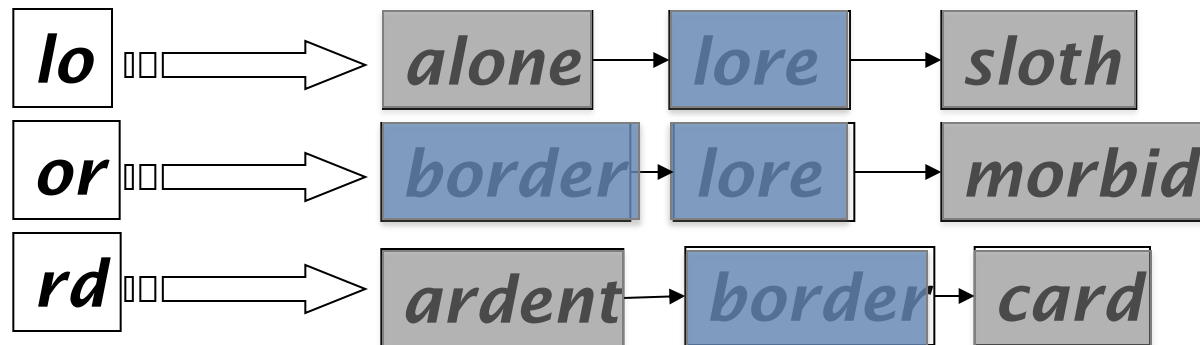
- A commonly-used measure of overlap
- Let X and Y be two sets; then the J.C. is

$$|X \cap Y| / |X \cup Y|$$

- Equals 1 when X and Y have the same elements and zero when they are disjoint
- X and Y don't have to be of the same size
- Always assigns a number between 0 and 1
 - Now threshold to decide if you have a match
 - E.g., if J.C. > 0.8, declare a match

Matching trigrams

- Consider the query *lord* – we wish to identify words matching 2 of its 3 bigrams (*lo*, *or*, *rd*)



Standard postings “merge” will enumerate ...

Adapt this to using Jaccard (or another) measure.

Context-sensitive spell correction

- Consider the phrase query “*flew form Heathrow*”
- We’d like the IR system to respond
Did you mean “*flew from Heathrow*”?
because no docs matched the query phrase.
- Need surrounding context to catch this error

Context-sensitive correction

- Query: *flew form heathrow*
 - We do not know which word(s) is an error.
- First idea: retrieve dictionary terms **close** (in weighted edit distance) to each query term
- Now try all possible resulting phrases with one word “fixed” at a time
 - *flew from heathrow* ✓✓ → 10K
 - *fled form heathrow* → 5K
 - *flea form heathrow* → 3K
 - ... and so on
- **Hit-based spelling correction**: Suggest the alternative that has lots of hits.

Exercise

- Suppose that for ***“flew form Heathrow”*** we have 7 alternatives for flew, 19 for form and 3 for heathrow. How many “corrected” phrases will we enumerate in this scheme?

General issues in spell correction

- We enumerate multiple alternatives for “Did you mean?”
- Need to figure out which to present to the user
 - The alternative hitting most docs
 - Query log analysis
- More generally, rank alternatives probabilistically

$$\operatorname{argmax}_{corr} P(corr \mid query)$$

- From Bayes rule, this is equivalent to

$$\operatorname{argmax}_{corr} P(query \mid corr) * P(corr)$$

Noisy channel

Language model

$P(\text{did you mean})$
 $>$
 $P(\text{mean you did})$

SOUNDEX

Soundex

- Class of heuristics to expand a query into **phonetic** equivalents
 - Language specific – mainly for names :
 - E.g., *chebyshev* → *tchebycheff*
- Invented for the U.S. census ... in 1918

Soundex – typical algorithm

- Turn every token to be indexed into a 4-character reduced form
- Do the same with query terms
- Build and search an index on the reduced forms
 - (when the query calls for a soundex match)
- <http://www.creativyst.com/Doc/Articles/SoundEx1/SoundEx1.htm#Top>

Soundex – typical algorithm

1. Retain the first letter of the word.
2. Change all occurrences of the following letters to '0' (zero):
'A', 'E', 'I', 'O', 'U', 'H', 'W', 'Y'.
3. Change letters to digits as follows:
 - B, F, P, V \rightarrow 1
 - C, G, J, K, Q, S, X, Z \rightarrow 2
 - D, T \rightarrow 3
 - L \rightarrow 4
 - M, N \rightarrow 5
 - R \rightarrow 6

Soundex continued

4. Remove all pairs of consecutive digits.
5. Remove all zeros from the resulting string.
6. Pad the resulting string with trailing zeros and return the first four positions, which will be of the form <uppercase letter> <digit> <digit> <digit>.

E.g., *Herman* becomes H655.

Will *hermann* generate the same code?

Soundex

- Soundex is the classic algorithm, provided by most databases (Oracle, Microsoft, ...)
- How useful is soundex?
- Not very – for information retrieval
- Okay for “high recall” tasks (e.g., Interpol), though biased to names of certain nationalities
- Zobel and Dart (1996) show that other algorithms for phonetic matching perform much better in the context of IR

What queries can we process?

- We have
 - Positional inverted index with skip pointers
 - Wild-card index
 - Spell-correction
 - Soundex
- Queries such as
(SPELL(moriset) /3 toron*to) OR SOUNDEX(chaikofski)