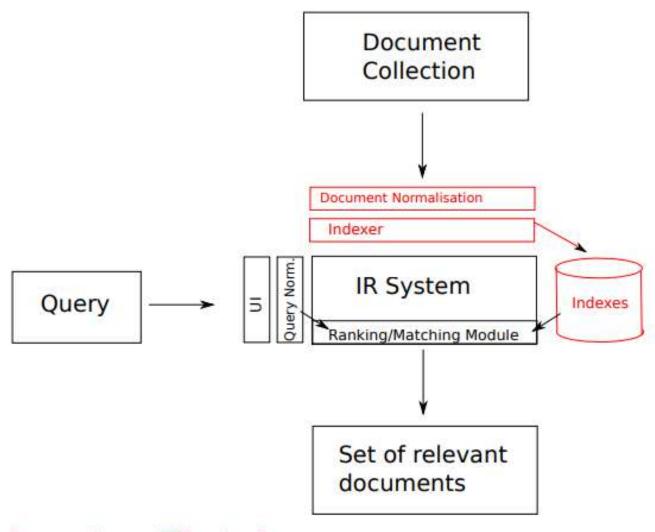
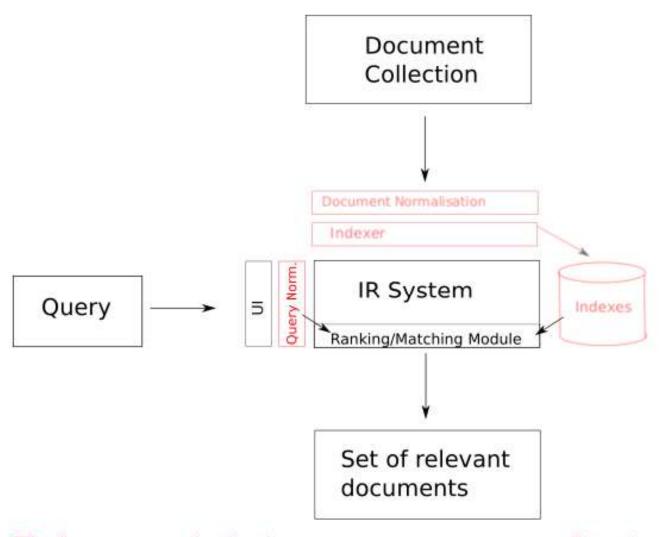
Lecture 5 Tolerant Retrieval

IR System components



Last time: The indexer

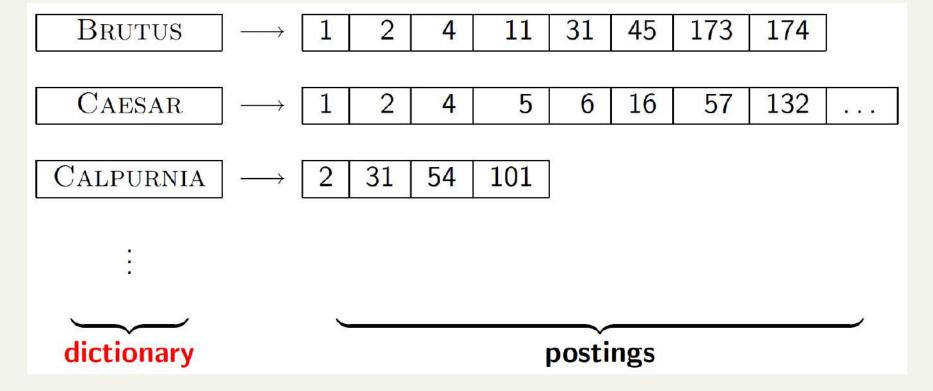
IR System components



Today: more indexing, some query normalisation

Dictionary data structures for inverted indexes

The dictionary data structure stores the term vocabulary, document frequency, pointers to each postings list ... in what data structure?



A naïve dictionary

An array of struct:

term	document	pointer to
	frequency	postings list
а	656,265	\longrightarrow
aachen	65	\longrightarrow
	***	3.4.3
zulu	221	\longrightarrow

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

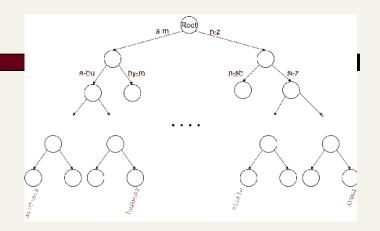
- Two main choices:
 - Hash table
 - Tree
- Some IR systems use hashes, some trees

Hashes

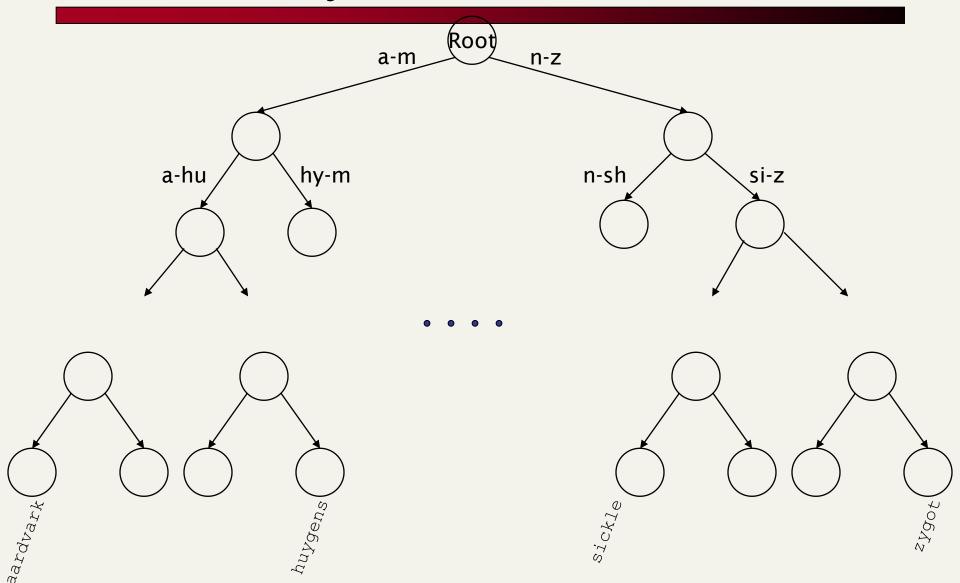
- Each vocabulary term is hashed to an integer
 - (We assume you've seen hashtables before)
- Pros:
 - Lookup is faster than for a tree: O(1)
- Cons:
 - No easy way to find minor variants:
 - judgment/judgement
 - No prefix search [tolerant retrieval]
 - If vocabulary keeps going, need to occasionally do the expensive operation of rehashing everything

Trees

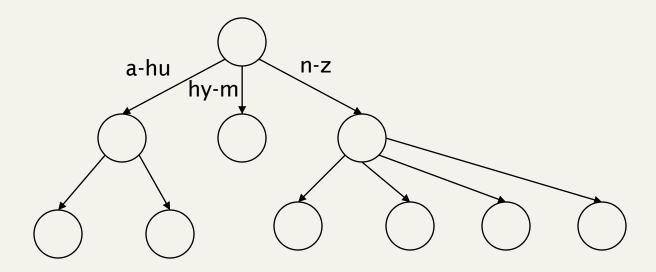
- Simplest: Binary tree
- More usual: B-trees
- Pros:
 - Solves the prefix problem (terms starting with hyp)
- Cons:
 - Slower: O(log M) [and this requires balanced tree]
 - Rebalancing binary trees is expensive
 - But B-trees mitigate the rebalancing problem



Tree: binary tree

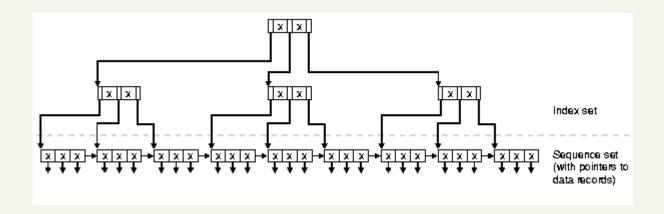


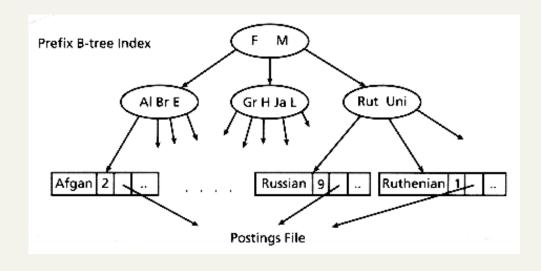
Tree: B-tree

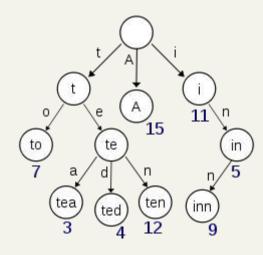


Definition: Every internal nodel has a number of children in the interval [a,b] where a, b are appropriate natural numbers.

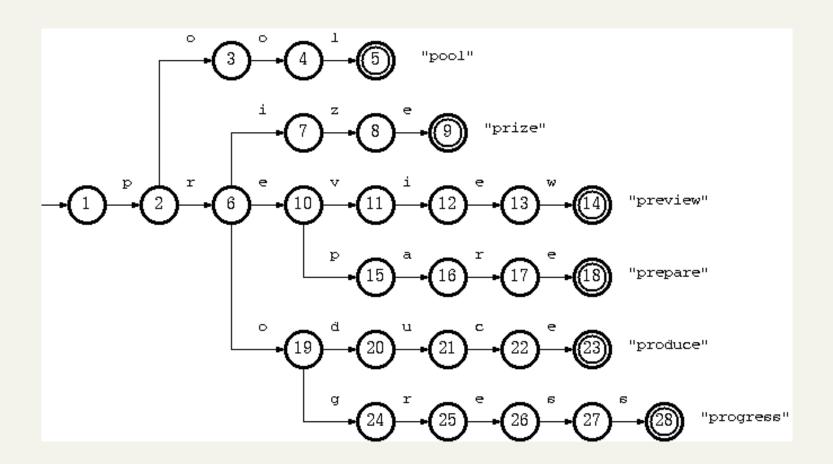
Tree: Dictionary







Tree: Dictionary



Keyword-based query language

- Single-word query
 - 1-word matching
 - Frequency of matching
- Content query
 - Given word near other words
 - Phrase
 - Proximity (near by)

Keyword-based query language

- Boolean Query
 - AND, OR, NOT
 - Query syntax tree
 - AND (satisfy both)
 - OR (satisfy either one)
 - A BUT B (satisfy A but not B)
 - Reasonable approach → "Fuzzy" technique

Keyword-based query language

Feedback technology

```
User → "query" → "result"

→ pick a doc → as a "query"

→ "results" ......
```

Wild-card queries

Wild-card queries: *

- mon*: find all docs containing any word beginning "mon".
- Easy with binary tree (or B-tree) lexicon: retrieve all words in range: mon ≤ w < moo</p>
- *mon: find words ending in "mon": harder
 - Maintain an additional B-tree for terms backwards.

Can retrieve all words in range: *nom* ≤ *w* < *non*.

Exercise: from this, how can we enumerate all terms meeting the wild-card query *pro*cent*?

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.
- E.g., consider the query:
 - se*ate AND fil*er

This may result in the execution of many Boolean *AND* queries.

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

- For term *hello*, index under:
 - hello\$, ello\$h, llo\$he, lo\$hel, o\$hell
 where \$ is a special symbol.
- 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*
```

Permuterm query processing

- Rotate query wild-card to the right
- Now use B-tree lookup as before.
- Permuterm problem: ≈ quadruples lexicon size

Empirical observation for English.

Bigram (k-gram) indexes

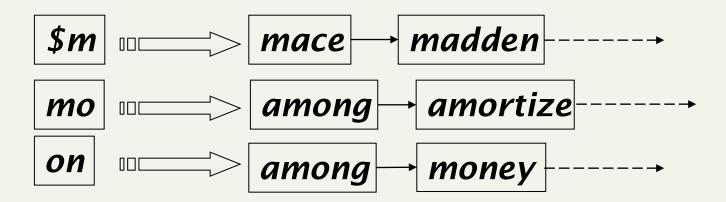
- Enumerate all k-grams (sequence of k chars)
 occurring in any term
- e.g., from text "April is the cruelest month" we get the 2-grams (bigrams)

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

- \$ is a special word boundary symbol
- Maintain a <u>second</u> inverted index <u>from bigrams to</u> <u>dictionary terms</u> that match each bigram.

Bigram index example

The k-gram index finds terms based on a query consisting of k-grams

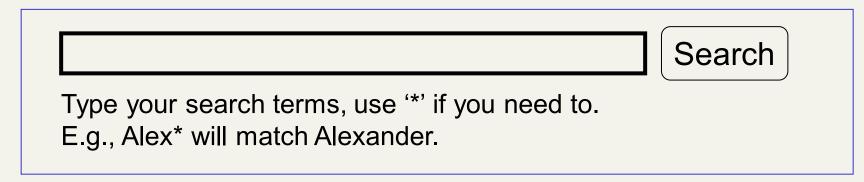


Processing *n*-gram wild-cards

- 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!



Does Google allow wildcard queries?

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
 - Will not catch typos resulting in correctly spelled words
 - e.g., $from \rightarrow form$
 - Context-sensitive
 - Look at surrounding words,
 - e.g., I flew form Heathrow to Narita.

Document correction

- Especially 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 (adjacent on the QWERTY keyboard, so more likely interchanged in typing).
- But also: web pages and even printed material has typos
- Goal: the dictionary contains fewer misspellings
- But often we don't change the documents but aim to fix the query-document mapping

Query mis-spellings

- Our principal focus here
 - E.g., the query *Alanis Morisett*
- 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 spelling correction

- There is a list of "correct" words for instance a standard dictionary (Webster's, OED...)
- Then we need a way of computing the distance between a misspelled word and a correct word
 - for instance Edit/Levenshtein distance
 - k-gram overlap
- Return the "correct" word that has the smallest distance to the misspelled word.

information → information

Misspelled Word Correction

When?

If a query word (combination) is quite rare or not available at all in the dictionary

Approach:

- Find similar term(s)
- Calculate their similarity to the query term
- 3. Choose the most frequent ones

Misspelled 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)

Misspelled 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.
- See http://www.merriampark.com/ld.htm for a nice example plus an applet.

Weighted edit distance

- As above, but the weight of an operation depends on the character(s) involved
 - Meant to capture OCR or keyboard errors, e.g. 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

Using edit distances

- 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

Edit distance to all dictionary terms?

- Given a (mis-spelled) query do we compute its edit distance to every dictionary term?
 - Expensive and slow
 - Alternative?
- How do we cut the set of candidate dictionary terms?
- One possibility is to use n-gram overlap for this
- This can also be used by itself for spelling correction.

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
- 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

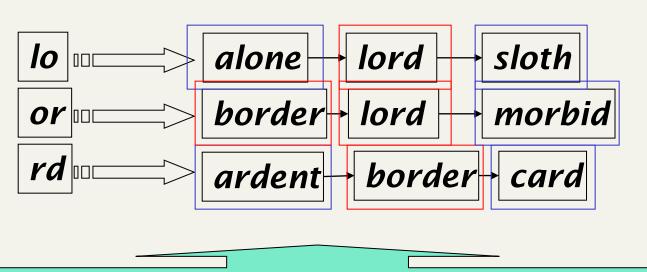
- A commonly-used measure of overlap
- Let X and Y be two sets; then the J.C. is

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

- 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

- Text: I flew from Heathrow to Narita.
- Consider the phrase query "flew form Heathrow"
- We'd like to respond

Did you mean "flew from Heathrow"?

because no docs matched the query phrase.

Context-sensitive correction

- Need surrounding context to catch this.
- 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
 - fled form heathrow
 - flea form heathrow
- 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?

Another approach

- Break phrase query into a conjunction of biwords (Lecture 2).
- Look for biwords that need only one term corrected.
- Enumerate phrase matches and ... rank them!

General issues in spell correction

- We enumerate multiple alternatives for "Did you mean?"
- Need to figure out which to present to the user
- Use heuristics
 - The alternative hitting most docs
 - Query log analysis + tweaking
 - For especially popular, topical queries
- Spell-correction is computationally expensive
 - Avoid running routinely on every query?
 - Run only on queries that matched few docs

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)

Soundex – typical algorithm

- 1. Retain the first letter of the word.
- 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 → 1
- $\bullet \quad \mathsf{C}, \, \mathsf{G}, \, \mathsf{J}, \, \mathsf{K}, \, \mathsf{Q}, \, \mathsf{S}, \, \mathsf{X}, \, \mathsf{Z} \to \mathsf{2}$
- $D,T \rightarrow 3$
- $L \rightarrow 4$
- \blacksquare M, N \rightarrow 5
- \blacksquare 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

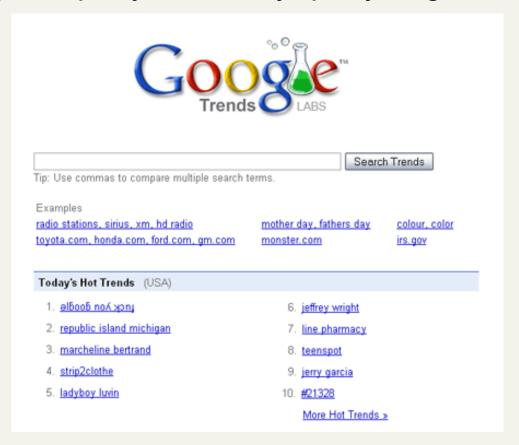
What queries can we process?

- We have
 - Positional inverted index with skip pointers
 - Wild-card index
 - Spell-correction
 - Soundex
- Queries such as

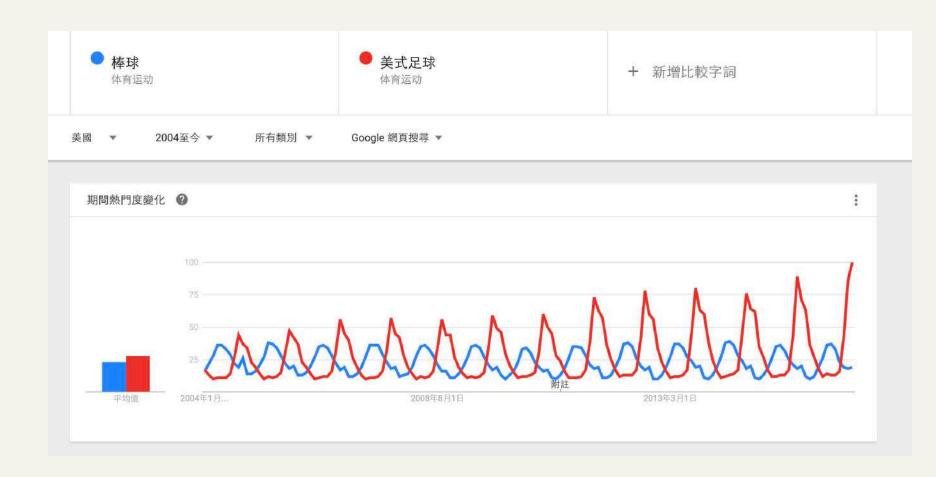
(SPELL(moriset) /3 toron*to) OR SOUNDEX(chaikofski)

About Query Logs

- Google Trend
 - logs of query issued by query engine



Google Trends



Why each country has each surveillance system?

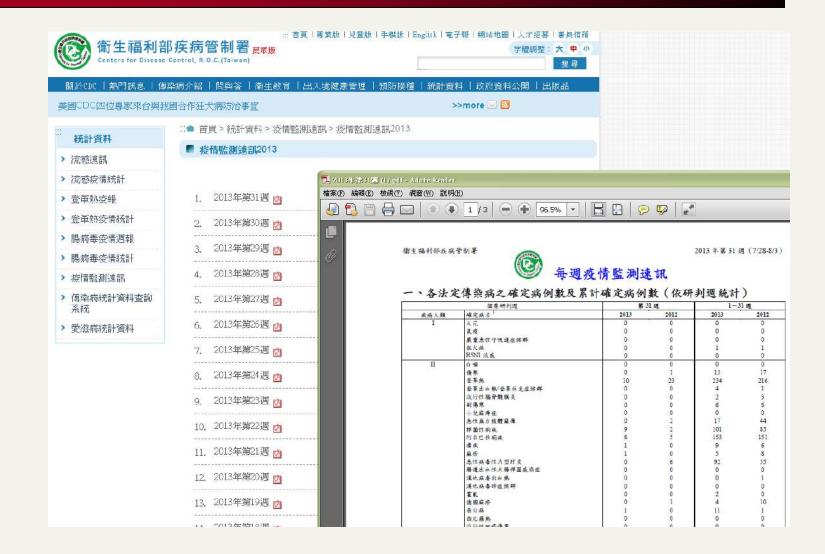
- Influenza epidemics are a major public health concern, because it causes tens of millions of illnesses each year.
- To reduce the victims, the early detection of influenza epidemics is a national mission in every country.
- BUT: These surveillance systems basically rely on hospital reports (written manually).



Centers for Disease Control and Prevention (CDC)



CDC Taiwan



Recent Approach

- using Phone Call data
 - Espino et al. (2003) used data of a telephone triage service, a public service, to give an advice to users via telephone. They reported the number of telephone calls that correlates with influenza epidemics.
- using Drug sale data
 - Magruder (2003) used the amount drug sales.

Among various approaches...

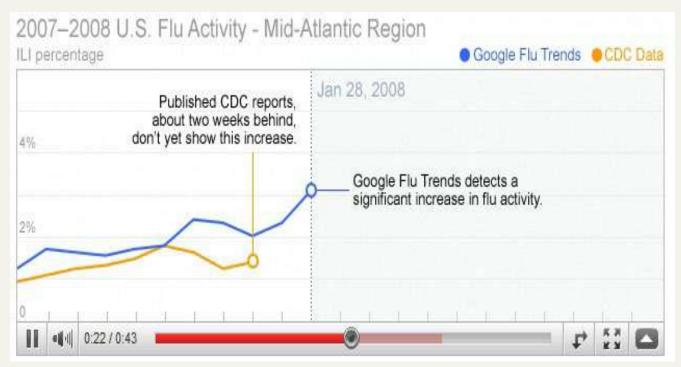
The State-of-the-Art Web based Approach

 Ginsberg et al. (Nature 2009) used Google web search queries that correlate with an influenza epidemic, such as "flu", "fever".



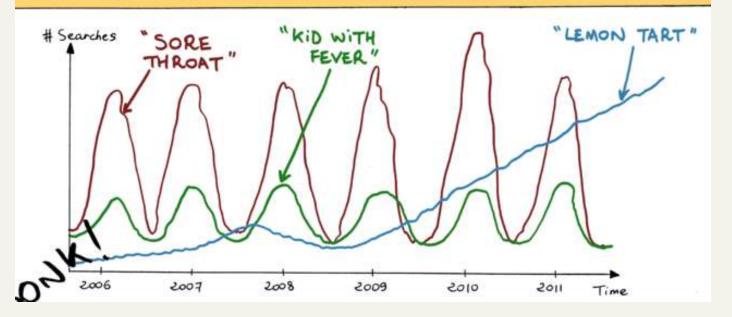


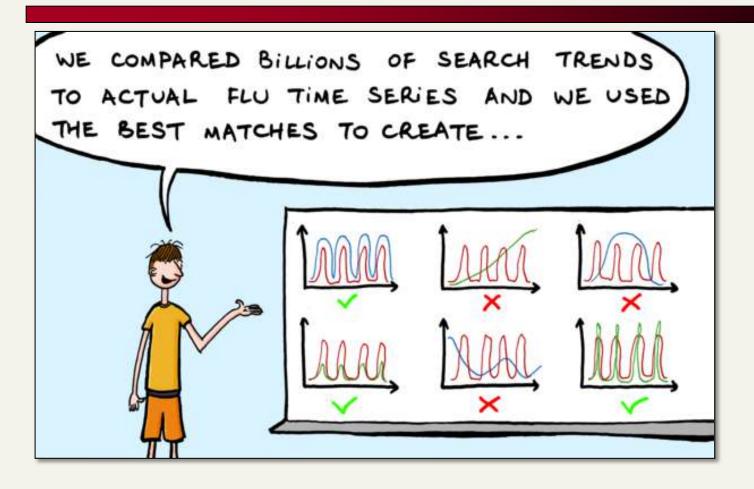
CDC vs Google Flu Trends?

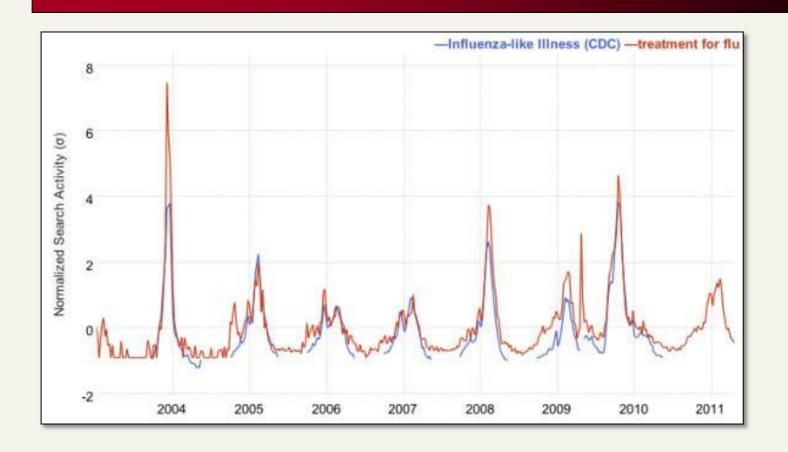


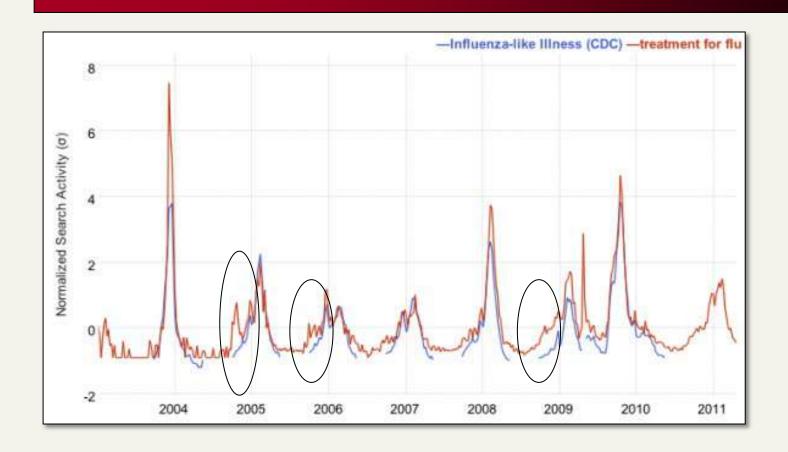
Source: http://www.google.org/flutrends/

BEING HUGE DATA NERDS, WE WONDERED IF TRENDS IN FLU-RELATED SEARCH ACTIVITY ON GOOGLE MIGHT REFLECT THE PATTERN OF ACTUAL FLU ACTIVITY.

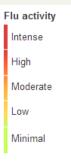








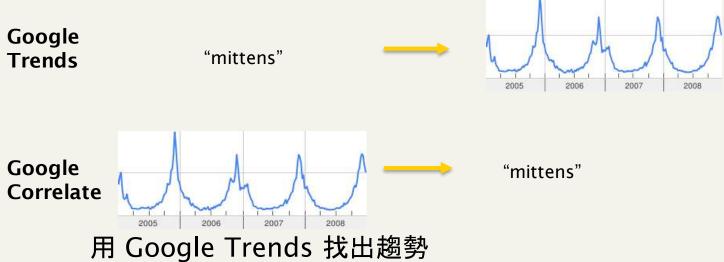




http://www.google.org/flutrends/

What is Google Correlate?

- With Google Trends, you type in a query and get back a series of its frequency.
- Google Correlate is like Google Trends in reverse.
- With Google Correlate, you enter a data series (the target) and get back queries whose frequency follows a similar pattern.



用 Google Trends 找出趨勢 再透過Google Correlate 找出與此趨勢相關的事物

Google Correlate

