# Introduction to
# **Information Retrieval**

http://nlp.stanford.edu/IR-book/

# Unstructured data in 1680

- Which plays of Shakespeare contain the words ***Brutus*** *AND* ***Caesar*** but *NOT* ***Calpurnia***?

- One could `grep` all of Shakespeare's plays for ***Brutus*** and ***Caesar,*** then strip out lines containing ***Calpurnia***?

- Why is that not the answer?
  - Slow (for large corpora)
  - <u>*NOT*</u> ***Calpurnia*** is non-trivial
  - Other operations (e.g., find the word ***Romans*** near ***countrymen***) not feasible
  - Ranked retrieval (best documents to return)
    - Later lectures

# Term-document incidence

|  | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 1 | 1 | 0 | 0 | 0 | 1 |
| **Brutus** | 1 | 1 | 0 | 1 | 0 | 0 |
| **Caesar** | 1 | 1 | 0 | 1 | 1 | 1 |
| **Calpurnia** | 0 | 1 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 1 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1 | 0 | 1 | 1 | 1 | 1 |
| **worser** | 1 | 0 | 1 | 1 | 1 | 0 |

***Brutus*** *AND* ***Caesar*** *BUT NOT* ***Calpurnia***

1 if play contains word, 0 otherwise

# Incidence vectors

- So we have a 0/1 vector for each term.

- To answer query: take the vectors for **Brutus, Caesar** and **Calpurnia** (complemented) ➔ bitwise *AND*.

- 110100 *AND* 110111 *AND* 101111 = 100100.

# Answers to query

- ## Antony and Cleopatra, Act III, Scene ii

*Agrippa* [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,

When Antony found Julius **Caesar** dead,

He cried almost to roaring; and he wept

When at Philippi he found **Brutus** slain.

- ## Hamlet, Act III, Scene ii

*Lord Polonius:* I did enact Julius **Caesar** I was killed i' the

Capitol; **Brutus** killed me.

# Basic assumptions of Information Retrieval

- Collection: Fixed set of documents

- Goal: Retrieve documents with information that is <u>relevant</u> to the user's information need and helps the user complete a task

# How good are the retrieved docs?

- *Precision* : Fraction of retrieved docs that are relevant to user's information need
- *Recall* : Fraction of relevant docs in collection that are retrieved

# Bigger collections

- Consider $N$ = 1 million documents, each with about 1000 words.

- Avg 6 bytes/word including spaces/punctuation

  - 6GB of data in the documents.

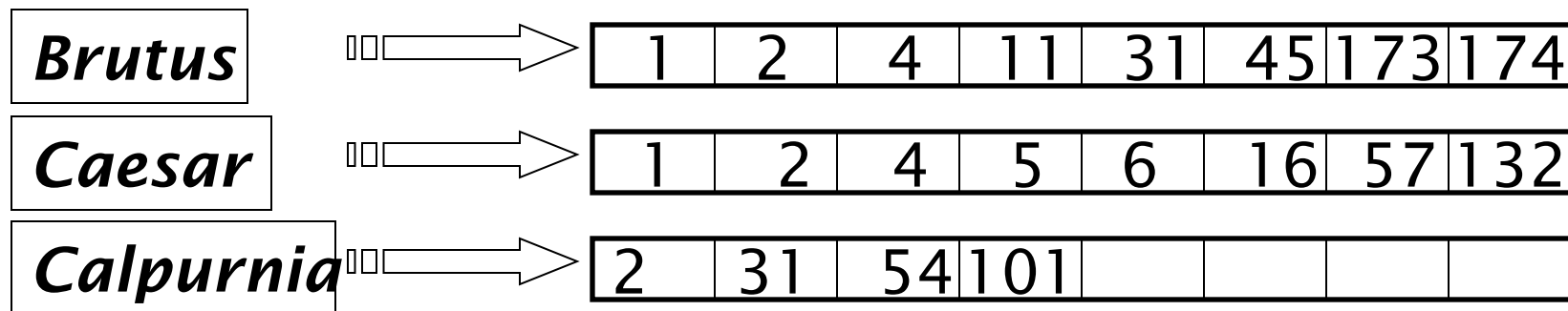- Say there are $M$ = 500K *distinct* terms among these.

# Can't build the matrix

- 500K x 1M matrix has half-a-trillion 0's and 1's.

- But it has no more than one billion 1's.   ⟵ Why?

  - matrix is extremely sparse.

- What's a better representation?

  - We only record the 1 positions.
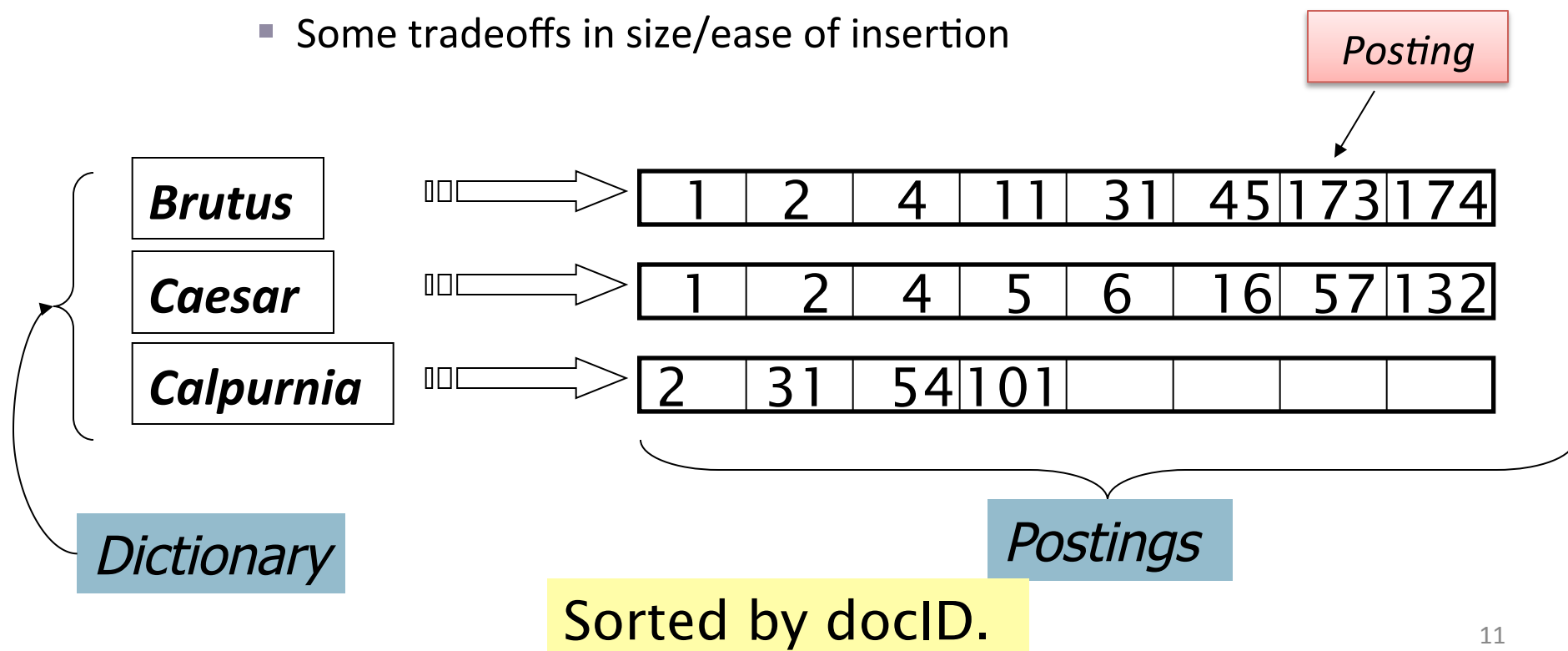
# Inverted index

- For each term *t*, we must store a list of all documents that contain *t*.

  - Identify each by a **docID**, a document serial number

- Can we use fixed-size arrays for this?

| **Brutus** | → | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|---|---|---|---|---|---|---|---|---|---|

| **Caesar** | → | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 |
|---|---|---|---|---|---|---|---|---|---|

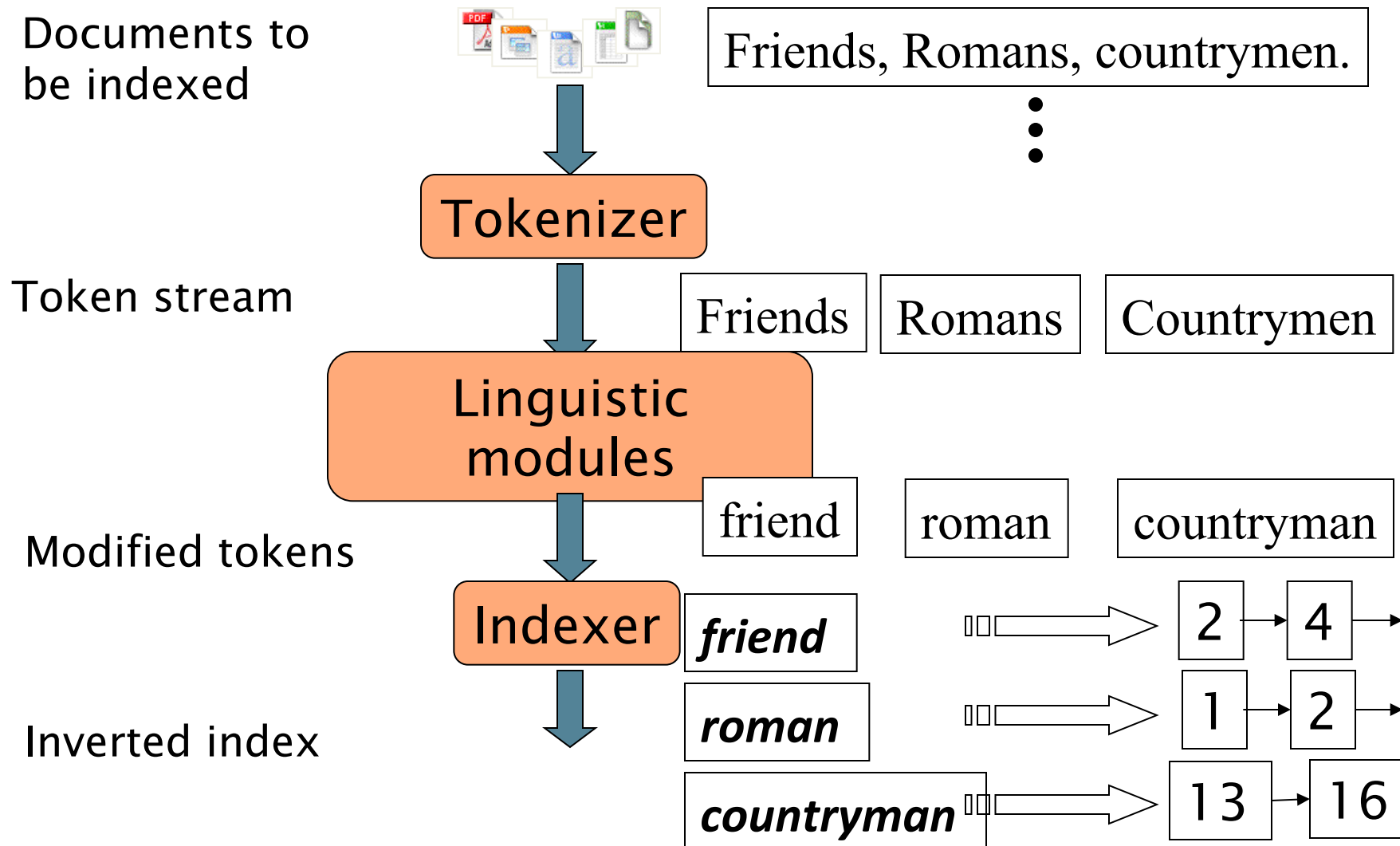| **Calpurnia** | → | 2 | 31 | 54 | 101 | | | | |
|---|---|---|---|---|---|---|---|---|---|

What happens if the word ***Caesar*** is added to document 14?

# Inverted index

- ## We need variable-size postings lists

  - ### On disk, a continuous run of postings is normal and best

  - ### In memory, can use linked lists or variable length arrays

    - #### Some tradeoffs in size/ease of insertion

*Posting*

| **Brutus** | | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |

| **Caesar** | | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 |

| **Calpurnia** | | 2 | 31 | 54 | 101 | | | | |

*Dictionary*

*Postings*

Sorted by docID.

11

# Inverted index construction

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

# Indexer steps: Token sequence

- Sequence of (Modified token, Document ID) pairs.

| Term | docID |
|------|------:|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |

Doc 1

Doc 2

I did enact Julius Caesar I was killed i' the Capitol; Brutus killed me.

So let it be with Caesar. The noble Brutus hath told you Caesar was ambitious

# Indexer steps: Sort

- ## Sort by terms
  - And then docID

**Core indexing step**

| Term | docID |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |

| Term | docID |
|------|-------|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |

# Indexer steps: Dictionary & Postings

- **Multiple term entries in a single document are merged.**

- **Split into Dictionary and Postings**
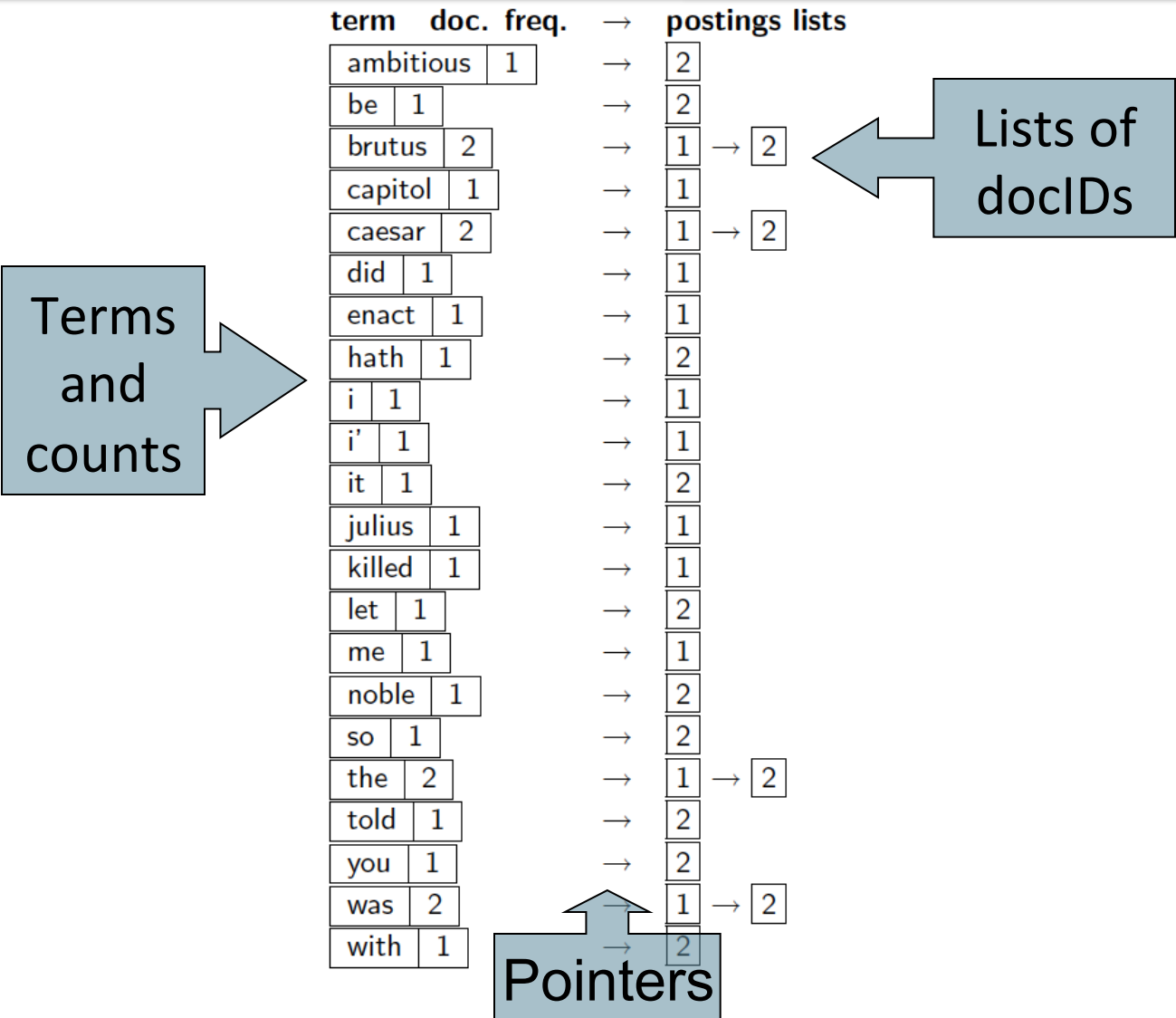
- **Doc. frequency information is added.**

  Why frequency? Will discuss later.

| Term | docID |
|---|---|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |

| term | doc. freq. | → | postings lists |
|---|---|---|---|
| ambitious | 1 | → | 2 |
| be | 1 | → | 2 |
| brutus | 2 | → | 1 → 2 |
| capitol | 1 | → | 1 |
| caesar | 2 | → | 1 → 2 |
| did | 1 | → | 1 |
| enact | 1 | → | 1 |
| hath | 1 | → | 2 |
| i | 1 | → | 1 |
| i' | 1 | → | 1 |
| it | 1 | → | 2 |
| julius | 1 | → | 1 |
| killed | 1 | → | 1 |
| let | 1 | → | 2 |
| me | 1 | → | 1 |
| noble | 1 | → | 2 |
| so | 1 | → | 2 |
| the | 2 | → | 1 → 2 |
| told | 1 | → | 2 |
| you | 1 | → | 2 |
| was | 2 | → | 1 → 2 |
| with | 1 | → | 2 |

# Where do we pay in storage?

| term | doc. freq. | → | postings lists |
|---|---|---|---|
| ambitious | 1 | → | 2 |
| be | 1 | → | 2 |
| brutus | 2 | → | 1 → 2 |
| capitol | 1 | → | 1 |
| caesar | 2 | → | 1 → 2 |
| did | 1 | → | 1 |
| enact | 1 | → | 1 |
| hath | 1 | → | 2 |
| i | 1 | → | 1 |
| i' | 1 | → | 1 |
| it | 1 | → | 2 |
| julius | 1 | → | 1 |
| killed | 1 | → | 1 |
| let | 1 | → | 2 |
| me | 1 | → | 1 |
| noble | 1 | → | 2 |
| so | 1 | → | 2 |
| the | 2 | → | 1 → 2 |
| told | 1 | → | 2 |
| you | 1 | → | 2 |
| was | 2 | → | 1 → 2 |
| with | 1 | → | 2 |

Lists of docIDs

Terms and counts

Pointers

16

# The index we just built

- How do we process a query?
  - Later - what kinds of queries can we process?

# Query processing: AND

- Consider processing the query:

  ***Brutus*** *AND* ***Caesar***

  - Locate ***Brutus*** in the Dictionary;

    - Retrieve its postings.

  - Locate *Caesar* in the Dictionary;

    - Retrieve its postings.

  - "Merge" the two postings:

| 2 | 4 | 8 | 16 | 32 | 64 | 128 | *Brutus* |
|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 | *Caesar* |

# The merge

- Walk through the two postings simultaneously, in time linear in the total number of postings entries



If list lengths are *x* and *y*, merge takes O(*x+y*) operations.
Crucial: postings sorted by docID.

# Intersecting two postings lists (a "merge" algorithm)

$$\text{INTERSECT}(p_1, p_2)$$

| | |
|---|---|
| 1 | $answer \leftarrow \langle \, \rangle$ |
| 2 | **while** $p_1 \neq \text{NIL}$ and $p_2 \neq \text{NIL}$ |
| 3 | **do if** $docID(p_1) = docID(p_2)$ |
| 4 | **then** $\text{ADD}(answer, docID(p_1))$ |
| 5 | $p_1 \leftarrow next(p_1)$ |
| 6 | $p_2 \leftarrow next(p_2)$ |
| 7 | **else if** $docID(p_1) < docID(p_2)$ |
| 8 | **then** $p_1 \leftarrow next(p_1)$ |
| 9 | **else** $p_2 \leftarrow next(p_2)$ |
| 10 | **return** $answer$ |

# Boolean queries: Exact match

- The Boolean retrieval model is being able to ask a query that is a Boolean expression:
  - Boolean Queries use *AND, OR* and *NOT* to join query terms
    - Views each document as a <u>set</u> of words
    - Is precise: document matches condition or not.
  - Perhaps the simplest model to build an IR system on
- Primary commercial retrieval tool for 3 decades.
- Many search systems you still use are Boolean:
  - Email, library catalog, Mac OS X Spotlight

# Example: WestLaw    http://www.westlaw.com/

- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992)

- Tens of terabytes of data; 700,000 users

- Majority of users *still* use boolean queries

- Example query:
  - What is the statute of limitations in cases involving the federal tort claims act?
  - LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM
    - /3 = within 3 words, /S = in same sentence

# Example: WestLaw    http://www.westlaw.com/

- Another example query:
  - Requirements for disabled people to be able to access a workplace
  - disabl! /p access! /s work-site work-place (employment /3 place)
- Note that SPACE is disjunction, not conjunction!
- Long, precise queries; proximity operators; incrementally developed; not like web search
- Many professional searchers still like Boolean search
  - You know exactly what you are getting
- But that doesn't mean it actually works better....

# Boolean queries:
# More general merges

- Exercise: Adapt the merge for the queries:

  ***Brutus*** *AND NOT* ***Caesar***

  ***Brutus*** *OR NOT* ***Caesar***

Can we still run through the merge in time O($x+y$)?

What can we achieve?

# Merging

What about an arbitrary Boolean formula?

*(**Brutus** OR **Caesar**) AND NOT*

*(**Antony** OR **Cleopatra**)*

- Can we always merge in "linear" time?
  - Linear in what?
- Can we do better?

# Query optimization

- What is the best order for query processing?

- Consider a query that is an *AND* of *n* terms.

- For each of the *n* terms, get its postings, then *AND* them together.

| Brutus | | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |
|---|---|---|---|---|---|---|---|---|---|

| Caesar | | 1 | 2 | 3 | 5 | 8 | 16 | 21 | 34 |
|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | | 13 | 16 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

Query: **Brutus** *AND* **Calpurnia** *AND* **Caesar**

26

# Query optimization example

- Process in order of increasing freq:
  - *start with smallest set, then keep cutting further.*

This is why we kept document freq. in dictionary

| **Brutus** | | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |
|---|---|---|---|---|---|---|---|---|---|

| **Caesar** | | 1 | 2 | 3 | 5 | 8 | 16 | 21 | 34 |
|---|---|---|---|---|---|---|---|---|---|

| **Calpurnia** | | 13 | 16 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

Execute the query as (**Calpurnia** *AND* **Brutus)** *AND* **Caesar**.

# More general optimization

- e.g., *(madding OR crowd) AND (ignoble OR strife)*
- Get doc. freq.'s for all terms.
- Estimate the size of each *OR* by the sum of its doc. freq.'s (conservative).
- Process in increasing order of *OR* sizes.

# Exercise

- Recommend a query processing order for

*(tangerine OR trees) AND (marmalade OR skies) AND (kaleidoscope OR eyes)*

| Term | Freq |
|------|------|
| eyes | 213312 |
| kaleidoscope | 87009 |
| marmalade | 107913 |
| skies | 271658 |
| tangerine | 46653 |
| trees | 316812 |

# Query processing exercises

- **Exercise**: If the query is ***friends*** *AND* ***romans*** *AND (NOT* ***countrymen****),* how could we use the freq of ***countrymen***?

- **Exercise**: Extend the merge to an arbitrary Boolean query. Can we always guarantee execution in time linear in the total postings size?

- **Hint**: Begin with the case of a Boolean *formula* query where each term appears only once in the query.

# Exercise

- Try the search feature at
  [http://www.rhymezone.com/shakespeare/](http://www.rhymezone.com/shakespeare/)
- Write down five search features you think it could do better

# What's ahead in IR?
# Beyond term search

- What about phrases?

  - ***Stanford University***

- Proximity: Find ***Gates*** *NEAR* ***Microsoft***.

  - Need index to capture position information in docs.

- Zones in documents: Find documents with (*author =* ***Ullman***) *AND* (text contains ***automata***).

# Evidence accumulation

- 1 vs. 0 occurrence of a search term
  - 2 vs. 1 occurrence
  - 3 vs. 2 occurrences, etc.
  - Usually more seems better
- Need term frequency information in docs

# Ranking search results

- Boolean queries give inclusion or exclusion of docs.

- Often we want to rank/group results

    - Need to measure proximity from query to each doc.

    - Need to decide whether docs presented to user are singletons, or a group of docs covering various aspects of the query.

# IR vs. databases:
# Structured vs unstructured data

- Structured data tends to refer to information in "tables"

| Employee | Manager | Salary |
|----------|---------|--------|
| Smith    | Jones   | 50000  |
| Chang    | Smith   | 60000  |
| Ivy      | Smith   | 50000  |

Typically allows numerical range and exact match (for text) queries, e.g.,
*Salary < 60000 AND Manager = Smith.*

35

# Unstructured data

- Typically refers to free-form text

- Allows

  - Keyword queries including operators

  - More sophisticated "concept" queries, e.g.,

    - find all web pages dealing with *drug abuse*

- Classic model for searching text documents

# Semi-structured data

- In fact almost no data is "unstructured"

- E.g., this slide has distinctly identified zones such as the *Title* and *Bullets*

- Facilitates "semi-structured" search such as

  - *Title* contains <u>data</u> AND *Bullets* contain <u>search</u>

… to say nothing of linguistic structure

# More sophisticated semi-structured search

- *Title* is about <u>Object Oriented Programming</u> AND *Author* something like <u>stro*rup</u>

- where * is the wild-card operator

- Issues:
  - how do you process "about"?
  - how do you rank results?

- The focus of XML search

# Clustering, classification and ranking

- **Clustering:** Given a set of docs, group them into clusters based on their contents.

- **Classification:** Given a set of topics, plus a new doc *D*, decide which topic(s) *D* belongs to.

- **Ranking:** Can we learn how to best order a set of documents, e.g., a set of search results

# The web and its challenges

- Unusual and diverse documents
- Unusual and diverse users, queries, information needs
- Beyond terms, exploit ideas from social networks
  - link analysis, clickstreams …

- How do search engines work?
  And how can we make them better?

# More sophisticated *information* retrieval

- Cross-language information retrieval

- Question answering

- Summarization

- Text mining

- …

# TOKENS AND TERMS

# 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*
    - *lowercase*, *lower-case*, *lower case* ?
    - It can be effective to get the user to put in possible hyphens
  - ***San Francisco***: one token or two?
    - How do you decide it is one token?

# Tokenization: language issues

- French
  - ***L'ensemble*** → one token or two?
    - *L* ? *L'* ? *Le* ?
    - Want *l'ensemble* to match with *un ensemble*
      - Until at least 2003, it didn't on Google
        - Internationalization!

- German noun compounds are not segmented
  - ***Lebensversicherungsgesellschaftsangestellter***
  - 'life insurance company employee'
  - German retrieval systems benefit greatly from a **compound splitter** module
    - Can give a 15% performance boost for German

44

# Tokenization: language issues

- Chinese and Japanese have no spaces between words:

  - 莎拉波娃现在居住在美国东南部的佛罗里达。

  - Not always guaranteed a unique tokenization

- Further complicated in Japanese, with multiple alphabets intermingled

  - Dates/amounts in multiple formats

フォーチュン500社は情報不足のため時間あた$500K(約6,000万円)

| Katakana | Hiragana | Kanji | Romaji |

End-user can express query entirely in hiragana!

# 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.

- $\quad\quad\quad\quad\quad$ ← → ← → $\quad\quad\quad\quad\quad\quad\quad$ ← 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

# Stop words

- With a stop list, you exclude from the dictionary entirely the commonest words. Intuition:
  - They have little semantic content: *the, a, and, to, be*
  - There are a lot of them: ~30% of postings for top 30 words
- But the trend is away from doing this:
  - Good compression techniques 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"

# Normalization to terms

- We need to "normalize" words in indexed text as well as query words into the same form

  - We want to match ***U.S.A.*** and ***USA***

- Result is terms: a term is a (normalized) word type, which is an entry in our IR system dictionary

- We most commonly implicitly define equivalence classes of terms by, e.g.,

  - deleting periods to form a term

    - ***U.S.A., USA*** ➔ ***USA***

  - deleting hyphens to form a term

    - ***anti-discriminatory, antidiscriminatory*** ➔ ***antidiscriminatory***

# Thesauri and soundex

- Do we handle synonyms and homonyms?
  - E.g., by hand-constructed equivalence classes
    - *car* = *automobile*      *color* = *colour*
  - We can rewrite to form equivalence-class terms
    - When the document contains *automobile*, index it under *car-automobile* (and vice-versa)
  - Or we can expand a query
    - When the query contains *automobile*, look under *car* as well
- What about spelling mistakes?
  - One approach is soundex, which forms equivalence classes of words based on phonetic heuristics

# 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 |

# PHRASE QUERIES AND POSITIONAL INDEXES

# Phrase queries

- Want to be able to answer queries such as "***stanford university***" – as a phrase

- Thus the sentence "*I went to university at Stanford*" is not a match.

  - The concept of phrase queries has proven easily understood by users; one of the few "advanced search" ideas that works

  - Many more queries are *implicit phrase queries*

- For this, it no longer suffices to store only

  *<term* : *docs>* entries

# A first attempt: Biword indexes

- Index every consecutive pair of terms in the text as a phrase

- For example the text "Friends, Romans, Countrymen" would generate the biwords
  - ***friends romans***
  - ***romans countrymen***

- Each of these biwords is now a dictionary term

- Two-word phrase query-processing is now immediate.

# Longer phrase queries

- Longer phrases are processed as we did with wild-cards:

- **stanford university palo alto** can be broken into the Boolean query on biwords:

**stanford university** AND **university palo** AND **palo alto**

Without the docs, we cannot verify that the docs matching the above Boolean query do contain the phrase.

Can have false positives!

# Extended biwords

- Parse the indexed text and perform part-of-speech-tagging (POST).
- Bucket the terms into (say) Nouns (N) and articles/ prepositions (X).
- Call any string of terms of the form NX*N an <u>extended biword</u>.
  - Each such extended biword is now made a term in the dictionary.
- Example:  ***catcher in the rye***

  **N        X   X    N**

- Query processing: parse it into N's and X's
  - Segment query into enhanced biwords
  - Look up in index: ***catcher rye***

# Issues for biword indexes

- False positives, as noted before

- Index blowup due to bigger dictionary
  - Infeasible for more than biwords, big even for them


- Biword indexes are not the standard solution (for all biwords) but can be part of a compound strategy

# Solution 2: Positional indexes

- In the postings, store for each **term** the position(s) in which tokens of it appear:

  <**term**, number of docs containing **term**;

  *doc1*: position1, position2 … ;

  *doc2*: position1, position2 … ;

  etc.>

# Positional index example

*<be*: 993427;
*1*: 7, 18, 33, 72, 86, 231;
*2*: 3, 149;
*4*: 17, 191, 291, 430, 434;
*5*: 363, 367, …>

Which of docs 1,2,4,5 could contain "*to be or not to be*"?

- For phrase queries, we use a merge algorithm recursively at the document level

- But we now need to deal with more than just equality

59

# Processing a phrase query

- Extract inverted index entries for each distinct term: **to, be, or, not.**
- Merge their *doc:position* lists to enumerate all positions with "**to be or not to be**".

  - **to***:*

    - *2*:1,17,74,222,551; *4:8,16,190,429,433;* 7:13,23,191; …

  - **be***:*

    - *1*:17,19; *4:17,191,291,430,434;* 5:14,19,101; …

- Same general method for proximity searches

# Positional index size

- You can compress position values/offsets

- Nevertheless, a positional index expands postings storage *substantially*

- Nevertheless, a positional index is now standardly used because of the power and usefulness of phrase and proximity queries … whether used explicitly or implicitly in a ranking retrieval system.

# Ranked retrieval

- Thus far, our queries have all been Boolean.
  - Documents either match or don't.
- Good for expert users with precise understanding of their needs and the collection.
  - Also good for applications: Applications can easily consume 1000s of results.
- Not good for the majority of users.
  - Most users incapable of writing Boolean queries (or they are, but they think it's too much work).
  - Most users don't want to wade through 1000s of results.
    - This is particularly true of web search.

# Problem with Boolean search: feast or famine

- Boolean queries often result in either too few (=0) or too many (1000s) results.

- Query 1: "*standard user dlink 650*" → 200,000 hits

- Query 2: "*standard user dlink 650 no card found*": 0 hits

- It takes a lot of skill to come up with a query that produces a manageable number of hits.

  - AND gives too few; OR gives too many

# Ranked retrieval models

- Rather than a set of documents satisfying a query expression, in ranked retrieval, the system returns an ordering over the (top) documents in the collection for a query

- Free text queries: Rather than a query language of operators and expressions, the user's query is just one or more words in a human language

- In principle, there are two separate choices here, but in practice, ranked retrieval has normally been associated with free text queries and vice versa

# Feast or famine: not a problem in ranked retrieval

- When a system produces a ranked result set, large result sets are not an issue
  - Indeed, the size of the result set is not an issue
  - We just show the top $k$ ( ≈ 10) results
  - We don't overwhelm the user

  - Premise: the ranking algorithm works

# Scoring as the basis of ranked retrieval

- We wish to return in order the documents most likely to be useful to the searcher
- How can we rank-order the documents in the collection with respect to a query?
- Assign a score – say in [0, 1] – to each document
- This score measures how well document and query "match".

# Query-document matching scores

- We need a way of assigning a score to a query/ document pair

- Let's start with a one-term query

- If the query term does not occur in the document: score should be 0

- The more frequent the query term in the document, the higher the score (should be)

- We will look at a number of alternatives for this.

# Take 1: Jaccard coefficient

- A commonly used measure of overlap of two sets *A* and *B*

- jaccard*(A,B)* = |*A* ∩ *B*| / |*A* ∪ *B*|

- jaccard*(A,A)* = 1

- jaccard*(A,B)* = 0 if *A* ∩ *B* = 0

- *A* and *B* don't have to be the same size.

- Always assigns a number between 0 and 1.

# Jaccard coefficient: Scoring example

- What is the query-document match score that the Jaccard coefficient computes for each of the two documents below?

- <u>Query</u>: *ides of march*

- <u>Document</u> 1: *caesar died in march*

- <u>Document</u> 2: *the long march*

# Issues with Jaccard for scoring

- It doesn't consider *term frequency* (how many times a term occurs in a document)

- Rare terms in a collection are more informative than frequent terms. Jaccard doesn't consider this information

- We need a more sophisticated way of normalizing for length

- Later in this lecture, we'll use $|A \cap B| / \sqrt{|A \cup B|}$

- . . . instead of |A ∩ B|/|A ∪ B| (Jaccard) for length normalization.

# Recall: Binary term-document incidence matrix

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 1 | 1 | 0 | 0 | 0 | 1 |
| **Brutus** | 1 | 1 | 0 | 1 | 0 | 0 |
| **Caesar** | 1 | 1 | 0 | 1 | 1 | 1 |
| **Calpurnia** | 0 | 1 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 1 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1 | 0 | 1 | 1 | 1 | 1 |
| **worser** | 1 | 0 | 1 | 1 | 1 | 0 |

Each document is represented by a binary vector $\in \{0,1\}^{|V|}$

# Term-document count matrices

- Consider the number of occurrences of a term in a document:

  - Each document is a count vector in $\mathbb{N}^v$: a column below

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 157 | 73 | 0 | 0 | 0 | 0 |
| **Brutus** | 4 | 157 | 0 | 1 | 0 | 0 |
| **Caesar** | 232 | 227 | 0 | 2 | 1 | 1 |
| **Calpurnia** | 0 | 10 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 57 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 2 | 0 | 3 | 5 | 5 | 1 |
| **worser** | 2 | 0 | 1 | 1 | 1 | 0 |

# *Bag of words* model

- Vector representation doesn't consider the ordering of words in a document

- *John is quicker than Mary* and *Mary is quicker than John* have the same vectors

- This is called the <u>bag of words</u> model.

- In a sense, this is a step back: The positional index was able to distinguish these two documents.

- We will look at "recovering" positional information later in this course.

- For now: bag of words model

# Term frequency tf

- The term frequency $tf_{t,d}$ of term $t$ in document $d$ is defined as the number of times that $t$ occurs in $d$.

- We want to use tf when computing query-document match scores. But how?

- Raw term frequency is not what we want:

  - A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.

  - But not 10 times more relevant.

- Relevance does not increase proportionally with term frequency.

NB: frequency = count in IR

# Log-frequency weighting

- The log frequency weight of term t in d is

$$
w_{t,d} = \begin{cases} 1 + \log_{10} \mathrm{tf}_{t,d}, & \text{if } \mathrm{tf}_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}
$$

- $0 \to 0$, $1 \to 1$, $2 \to 1.3$, $10 \to 2$, $1000 \to 4$, etc.

- Score for a document-query pair: sum over terms $t$ in both $q$ and $d$:

- score $= \sum_{t \in q \cap d} (1 + \log \mathrm{tf}_{t,d})$

- The score is 0 if none of the query terms is present in the document.

# Document frequency

- Rare terms are more informative than frequent terms
  - Recall stop words

- Consider a term in the query that is rare in the collection (e.g., *arachnocentric*)

- A document containing this term is very likely to be relevant to the query *arachnocentric*

- → We want a high weight for rare terms like *arachnocentric*.

# Document frequency, continued

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., *high, increase, line*)
- A document containing such a term is more likely to be relevant than a document that doesn't
- But it's not a sure indicator of relevance.
- → For frequent terms, we want high positive weights for words like *high, increase, and line*
- But lower weights than for rare terms.
- We will use document frequency (df) to capture this.

# idf weight

- df$_t$ is the <u>document</u> frequency of *t*: the number of documents that contain *t*
  - df$_t$ is an inverse measure of the informativeness of *t*
  - df$_t$ ≤ *N*

- We define the idf (inverse document frequency) of *t* by

$$\text{idf}_t = \log_{10}(N/\text{df}_t)$$

  - We use log (*N*/df$_t$) instead of *N*/df$_t$ to "dampen" the effect of idf.

# idf example, suppose $N$ = 1 million

| term | $df_t$ | $idf_t$ |
|---|---:|---:|
| calpurnia | 1 | 6 |
| animal | 100 | 4 |
| sunday | 1,000 | 3 |
| fly | 10,000 | 2 |
| under | 100,000 | 1 |
| the | 1,000,000 | 0 |

$$\text{idf}_t = \log_{10}(N/\text{df}_t)$$

There is one idf value for each term $t$ in a collection.

# Effect of idf on ranking

- Does idf have an effect on ranking for one-term queries, like

  - iPhone

- idf has no effect on ranking one term queries

  - idf affects the ranking of documents for queries with at least two terms

  - For the query capricious person, idf weighting makes occurrences of capricious count for much more in the final document ranking than occurrences of person.

# Collection vs. Document frequency

- The collection frequency of *t* is the number of occurrences of *t* in the collection, counting multiple occurrences.

- Example:

| Word | Collection frequency | Document frequency |
|------|---------------------:|-------------------:|
| *insurance* | 10440 | 3997 |
| *try* | 10422 | 8760 |

- Which word is a better search term (and should get a higher weight)?

# tf-idf weighting

- The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$w_{t,d} = \log(1 + tf_{t,d}) \times \log_{10}(N / df_t)$$

- Best known weighting scheme in information retrieval
  - Note: the "-" in tf-idf is a hyphen, not a minus sign!
  - Alternative names: tf.idf, tf x idf
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection

# Score for a document given a query

$$\text{Score}(q,d) = \sum_{t \in q \cap d} \text{tf.idf}_{t,d}$$

- There are many variants
  - How "tf" is computed (with/without logs)
  - Whether the terms in the query are also weighted
  - …

# Binary → count → weight matrix

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 5.25 | 3.18 | 0 | 0 | 0 | 0.35 |
| **Brutus** | 1.21 | 6.1 | 0 | 1 | 0 | 0 |
| **Caesar** | 8.59 | 2.54 | 0 | 1.51 | 0.25 | 0 |
| **Calpurnia** | 0 | 1.54 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 2.85 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1.51 | 0 | 1.9 | 0.12 | 5.25 | 0.88 |
| **worser** | 1.37 | 0 | 0.11 | 4.15 | 0.25 | 1.95 |

Each document is now represented by a real-valued vector of tf-idf weights $\in \mathbb{R}^{|V|}$

# Documents as vectors

- So we have a |V|-dimensional vector space
- Terms are axes of the space
- Documents are points or vectors in this space
- Very high-dimensional: tens of millions of dimensions when you apply this to a web search engine
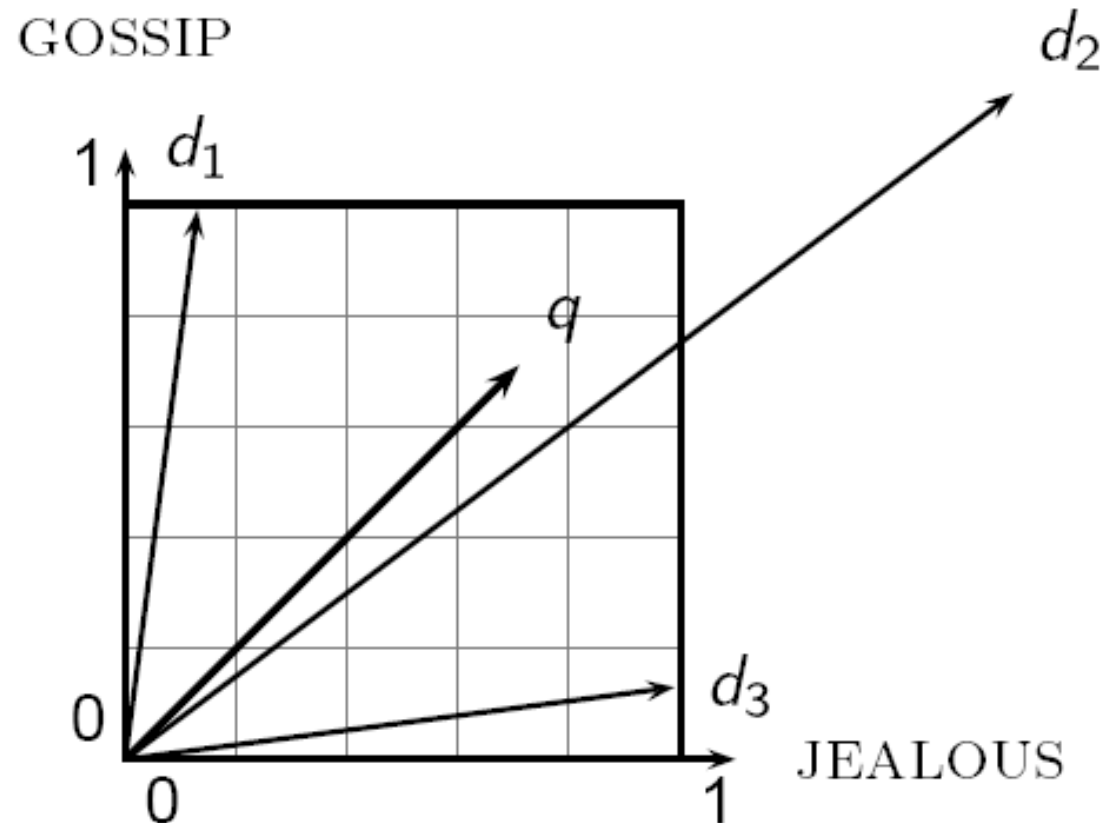- These are very sparse vectors - most entries are zero.

# Queries as vectors

- **Key idea 1:** Do the same for queries: represent them as vectors in the space

- **Key idea 2:** Rank documents according to their proximity to the query in this space

- proximity = similarity of vectors

- proximity ≈ inverse of distance

- Recall: We do this because we want to get away from the you're-either-in-or-out Boolean model.

- Instead: rank more relevant documents higher than less relevant documents

# Formalizing vector space proximity

- First cut: distance between two points
  - ( = distance between the end points of the two vectors)
- Euclidean distance?
- Euclidean distance is a bad idea . . .
- . . . because Euclidean distance is large for vectors of different lengths.

# Why distance is a bad idea

The Euclidean distance between $\vec{q}$ and $\vec{d_2}$ is large even though the

distribution of terms in the query $\vec{q}$ and the distribution of

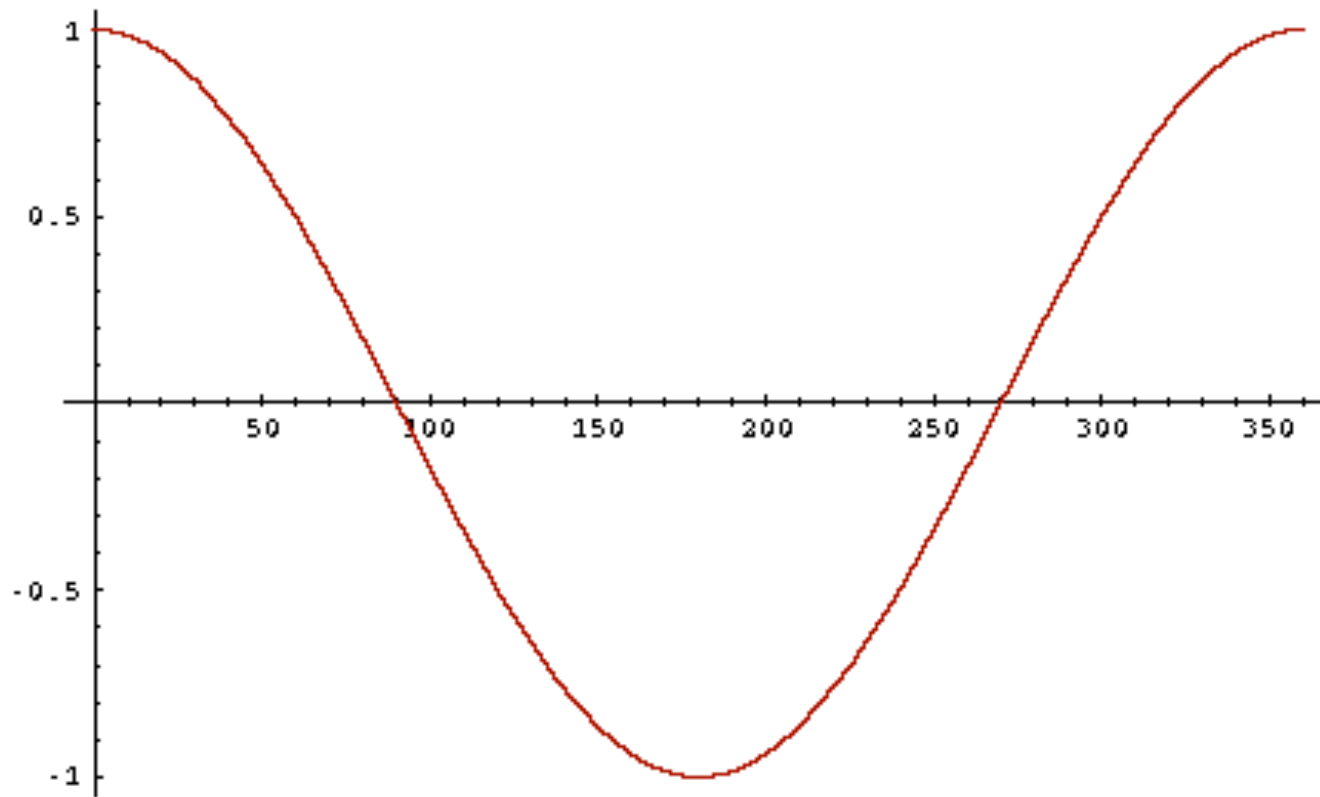terms in the document $\vec{d_2}$ are

very similar.

# Use angle instead of distance

- Thought experiment: take a document $d$ and append it to itself. Call this document $d$'.
- "Semantically" d and d' have the same content
- The Euclidean distance between the two documents can be quite large
- The angle between the two documents is 0, corresponding to maximal similarity.

- Key idea: Rank documents according to angle with query.

# From angles to cosines

- The following two notions are equivalent.
    - Rank documents in <u>decreasing</u> order of the angle between query and document
    - Rank documents in <u>increasing</u> order  of cosine (query,document)
- Cosine is a monotonically decreasing function for the interval [$0^o$, $180^o$]

# From angles to cosines



■ But how – *and why* – should we be computing cosines?

# Length normalization

- A vector can be (length-) normalized by dividing each of its components by its length – for this we use the $L_2$ norm:

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

- Dividing a vector by its $L_2$ norm makes it a unit (length) vector (on surface of unit hypersphere)

- Effect on the two documents d and d' (d appended to itself) from earlier slide: they have identical vectors after length-normalization.

  - Long and short documents now have comparable weights

# cosine(query,document)

Dot product

Unit vectors

$$\cos(\vec{q},\vec{d}) = \frac{\vec{q} \bullet \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\vec{q}}{|\vec{q}|} \bullet \frac{\vec{d}}{|\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

$q_i$ is the tf-idf weight of term *i* in the query
$d_i$ is the tf-idf weight of term *i* in the document

$\cos(\vec{q},\vec{d})$ is the cosine similarity of $\vec{q}$ and $\vec{d}$ … or, equivalently, the cosine of the angle between $\vec{q}$ and $\vec{d}$.
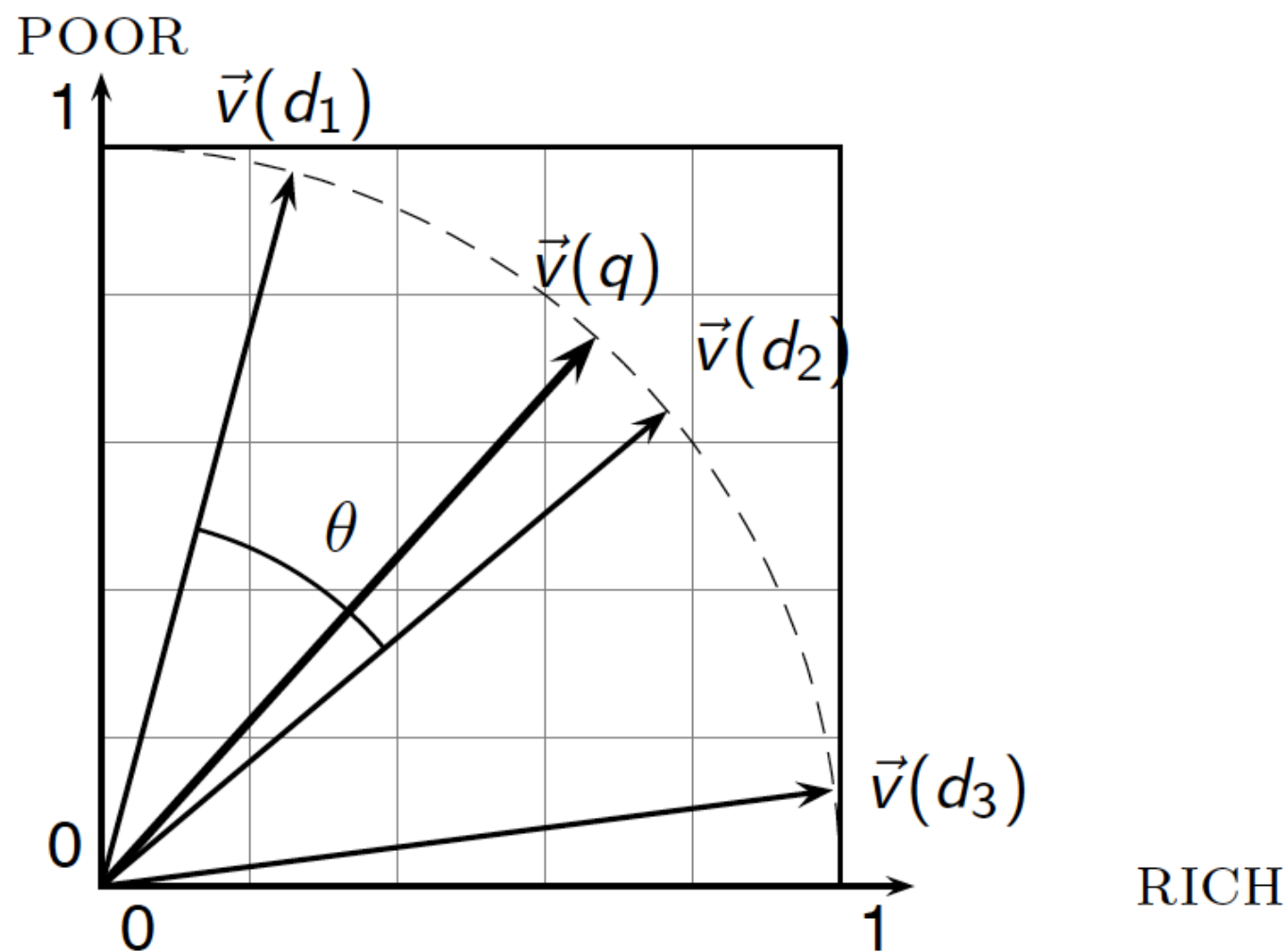
# Cosine for length-normalized vectors

- For length-normalized vectors, cosine similarity is simply the dot product (or scalar product):

$$\cos(\vec{q}, \vec{d}) = \vec{q} \bullet \vec{d} = \sum_{i=1}^{|V|} q_i d_i$$

for q, d length-normalized.

# Cosine similarity illustrated

# Cosine similarity amongst 3 documents

How similar are the novels

SaS: *Sense and Sensibility*

PaP: *Pride and Prejudice,* and

WH: *Wuthering Heights*?

| term | SaS | PaP | WH |
|------|-----|-----|-----|
| affection | 115 | 58 | 20 |
| jealous | 10 | 7 | 11 |
| gossip | 2 | 0 | 6 |
| wuthering | 0 | 0 | 38 |

## Term frequencies (counts)

Note: To simplify this example, we don't do idf weighting.

# 3 documents example contd.

**Log frequency weighting**

| term | SaS | PaP | WH |
|------|-----|-----|-----|
| affection | 3.06 | 2.76 | 2.30 |
| jealous | 2.00 | 1.85 | 2.04 |
| gossip | 1.30 | 0 | 1.78 |
| wuthering | 0 | 0 | 2.58 |

**After length normalization**

| term | SaS | PaP | WH |
|------|-----|-----|-----|
| affection | 0.789 | 0.832 | 0.524 |
| jealous | 0.515 | 0.555 | 0.465 |
| gossip | 0.335 | 0 | 0.405 |
| wuthering | 0 | 0 | 0.588 |

cos(SaS,PaP) ≈

$0.789 \times 0.832 + 0.515 \times 0.555 + 0.335 \times 0.0 + 0.0 \times 0.0$

≈ 0.94

cos(SaS,WH) ≈ 0.79

cos(PaP,WH) ≈ 0.69

Why do we have cos(SaS,PaP) > cos(SaS,WH)?