

CAIM: Cerca i Anàlisi d'Informació Massiva

FIB, Grau en Enginyeria Informàtica

Slides by Marta Arias, José Luis Balcázar,
Ramon Ferrer-i-Cancho, Ricard Gavaldá
Department of Computer Science, UPC

Fall 2018

<http://www.cs.upc.edu/~caim>

3. Implementation

Query answering

A **bad** algorithm:

```
input query  $q$ ;  
for every document  $d$  in database  
    check if  $d$  matches  $q$ ;  
    if so, add its docid to list  $L$ ;  
output list  $L$  (perhaps sorted in some way);
```

Query processing time should be largely independent of database size.

Probably proportional to answer size.

Central Data Structure

From terms to documents

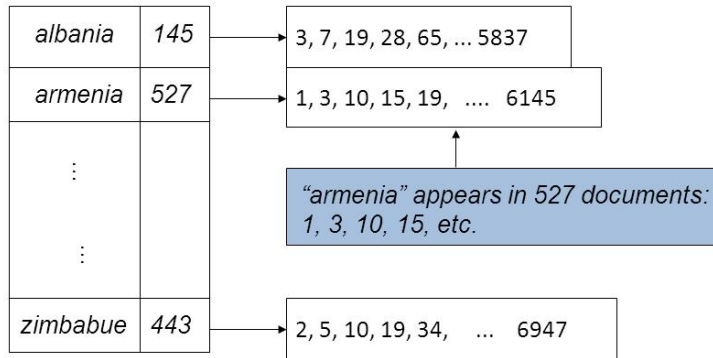
A **vocabulary** or **lexicon** or **dictionary**, usually kept in main memory, maintains all the indexed terms (*set*, *map*...); and, besides...

The Inverted File

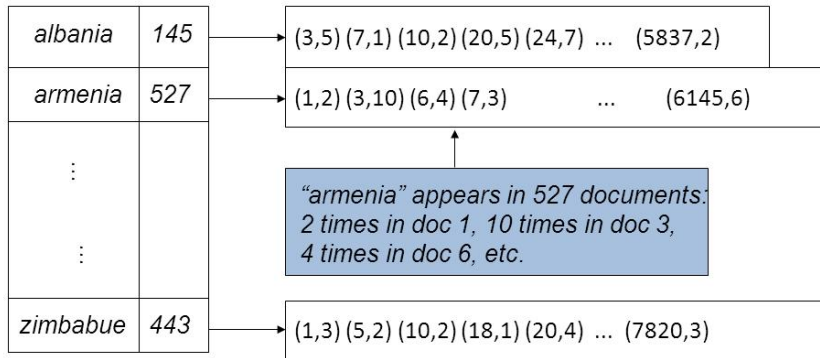
The crucial data structure for indexing.

- ▶ A data structure to support the operation:
 - ▶ “given term t , get all the documents that contain it”.
- ▶ The inverted file must support this operation (and variants) **very efficiently**.
- ▶ Built at preprocessing time, not at query time: can afford to spend some time in its construction.

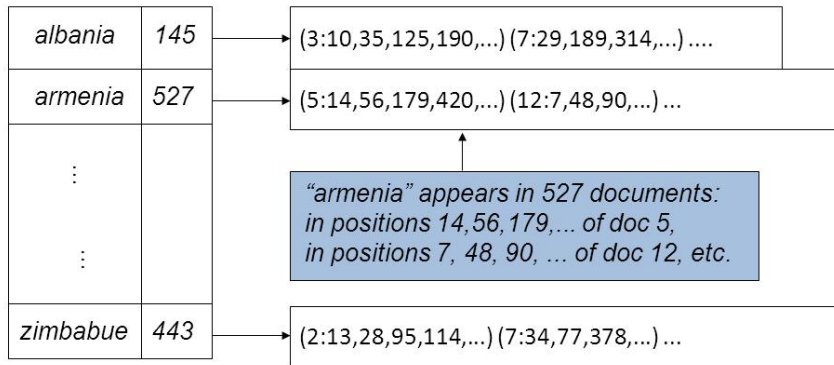
The inverted file: Variant 1



The inverted file: Variant 2



The inverted file: Variant 3



Postings

The inverted file is made of incidence/posting lists

We assign a *document identifier*, **docid** to each document.
The **dictionary** may fit in RAM for medium-size applications.

For each indexed term

a **posting list**: list of docid's (plus maybe other info) where the term appears.

- ▶ Wonderful if it fits in memory, but this is unlikely.
- ▶ Additionally: posting lists are
 - ▶ almost always sorted by *docid*
 - ▶ often compressed: minimize info to bring from disk!

Implementation of the Boolean Model, I

Simplest: Traverse posting lists

Conjunctive query: a AND b

- ▶ intersect the **posting lists** of a and b ;
- ▶ if sorted: can do a **merge-like intersection**;
- ▶ **time**: order of the **sum** of the lengths of posting lists.

intersect(input lists $L1$, $L2$, output list L):

```
while ( not  $L1.end()$  and not  $L2.end()$  )  
  if ( $L1.current() < L2.current()$ )  $L1.advance()$ ;  
  else if ( $L1.current() > L2.current()$ )  $L2.advance()$ ;  
  else {  $L.append(L1.current())$ ;  
         $L1.advance()$ ;  $L2.advance()$ ; }
```

Implementation of the Boolean Model, II

Simplest

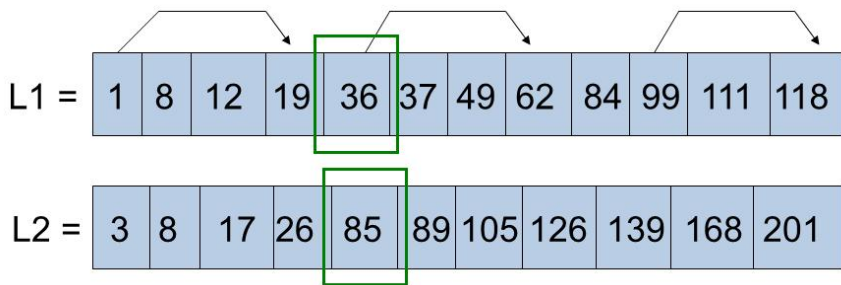
- ▶ Similar **merge-like union** for OR.
 - ▶ **Time**: again order of the **sum** of lengths of posting lists.
- ▶ Alternative: traverse one list and look up every docid in the other via **binary search**.
 - ▶ **Time**: length of shortest list **times** log of length of longest.

Example:

- ▶ $|L1| = 1000, |L2| = 1000$:
 - ▶ sequential scan: 2000 comparisons,
 - ▶ binary search: $1000 * 10 = 10,000$ comparisons.
- ▶ $|L1| = 100, |L2| = 10,000$:
 - ▶ sequential scan: 10,100 comparisons,
 - ▶ binary search: $100 * \log(10,000) = 1400$ comparisons.

Implementation of the Boolean Model, III

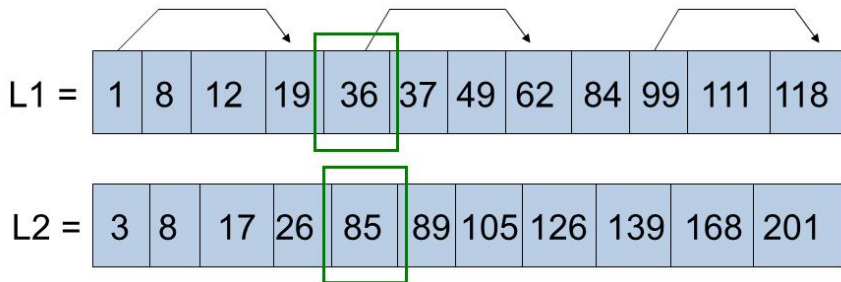
Sublinear time intersection: Skip pointers



- ▶ We've merged 1...19 and 3...26.
- ▶ We are looking at 36 and 85.
- ▶ Since $\text{pointer}(36)=62 < 85$, we can jump to 84 in L1.

Implementation of the Boolean Model, IV

Sublinear time intersection: Skip pointers



- ▶ Forward pointer from some elements.
- ▶ Either jump to next segment, or search within next segment (once).
- ▶ Optimal: in RAM, $\sqrt{|L|}$ pointers of length $\sqrt{|L|}$.
- ▶ Difficult to do well, particularly if the lists are **on disk**.

Query Optimization and Cost Estimation, I

Queries can be evaluated according to different plans

E.g. $a \text{ AND } b \text{ AND } c$ as

- ▶ $(a \text{ AND } b) \text{ AND } c$
- ▶ $(b \text{ AND } c) \text{ AND } a$
- ▶ $(a \text{ AND } c) \text{ AND } b$

E.g. $(a \text{ AND } b) \text{ OR } (a \text{ AND } c)$ also as

- ▶ $a \text{ AND } (b \text{ OR } c)$

The cost of an **execution plan** depends on the sizes of the lists **and** the sizes of intermediate lists.

Query Optimization and Cost Estimation, II

Example

Query: $(a \text{ AND } b) \text{ OR } (a \text{ AND } c \text{ AND } d)$.

Assume: $|La| = 3000$, $|Lb| = 1000$, $|Lc| = 2500$, $|Ld| = 300$.

- ▶ Three intersections plus one union, in the order given: up to cost 13600.
- ▶ Instead, $((d \text{ AND } c) \text{ AND } a)$: reduces to up to cost 11400.
- ▶ Rewrite to $a \text{ AND } (b \text{ OR } (c \text{ AND } d))$: reduces to up to cost 8400.

Implementation of the Vectorial Model, I

Problem statement

Fixed similarity measure $sim(d, q)$:

Retrieve

documents d_i which have a similarity to the query q

- ▶ either
 - ▶ above a threshold sim_{min} , or
 - ▶ the top r according to that similarity, or
 - ▶ all documents,
- ▶ sorted by decreasing similarity to the query q .

Must react **very fast** (thus, careful to the interplay with disk!), and with a reasonable memory expense.

Implementation of the Vectorial Model, II

Obvious nonsolution

Traverse all the documents, look at their terms in order to compute similarity, filter according to sim_{min} , and sort them...

... will not work.

Implementation of the Vectorial Model, III

Observations

Most documents include a **small proportion** of the available terms.

Queries usually include a **humanly small number** of terms.

Only a very **small proportion** of the documents will be relevant.

A priori bound r on the size of the answer known.

Inverted file available!

Implementation of the Vectorial Model, IV

Idea

Invert the loops:

- ▶ Outer loop on the terms t that appear in the query.
- ▶ Inner loop on documents that contain term t .
 - ▶ the reason for inverted index!
- ▶ Accumulate similarity for visited documents.
- ▶ Upon termination, normalize and sort.

Many additional subtleties can be incorporated.

Index compression, I

Why?

A large part of the query-answering time is spent

bringing posting lists from disks to RAM.

Need to minimize amount of bits to transfer.

Index compression schemes use:

- ▶ Docid's sorted in increasing order.
- ▶ Frequencies usually very small numbers.
- ▶ Can do better than e.g. 32 bits for each.

Index compression, II

Topic for self-study. At least:

- ▶ Unary self-delimiting code.
- ▶ Gap compression + Elias Gamma code.
- ▶ Continuation bit.
- ▶ Typical compression ratios.

E.g. books listed in the Presentation part of these notes.