Inverted Index

Lecture 12



Outline

- Background & Motivation
 - Full-Text Search
 - Big-O Review
 - Indexing
- The Inverted Index
 - An Example
 - Design a relational index
 - Advanced Issues
- Example in Cognitive Modeling



Derbinsky

Problem: Full-Text Search

- Given: set of "documents" containing "words"
 General problem in the field of *Information Retrieval*
- Task: find "best" document(s) that contain a set of words
- Requirements
 - Fast & scalable
 - Relevant results (precision, recall, f-score)
 - Expressive queries
 - Up-to-date



Example: Web Search





Other Examples





Scaling Up: # of "Documents"





1 December 2014

Scaling Up: Search Frequency





Big-O Review

- What does O(n) mean?
- A linear algorithm for full-text search?
 Find 5 documents that contain "WIT"
- What is the complexity in terms of documents (d) and average-words-perdocument (w)?



Is Linear-Time Google Possible?

• Assume simple query:

wentworth

ψQ

~4M Blu-ray

- Single listing of all web pages + words
 60T pages * 250 w/page * 8 bytes/w ~ 106PB
- Require 1s response
 - 7.5M * 2GHz 64-bit CPU (assume 1 cycle/w)
 - (7.5M * 68K) CPUs * 85W/CPU * \$.15/kWh ~ **\$1.8M/s**





Indexing

- Improve search speed at the cost of extra...
 - Memory for data structure(s)
 - Time to update the data structure(s)
- Backbone of databases (physical design)
 - Search engines
 - Graphics/game engines
 - Simulation software



. . .

Let's try some queries:

Inverted Index by Example

Given documents { D_1 , D_2 , D_3 }:

- $-D_1 =$ "it is what it is"
- $-D_2$ = "what is it"
- $-D_3 =$ "it is a banana"





Design a **Relational** Inverted Index

Develop a set of table(s) and index(es) that support efficient construction and querying of an inverted index

<u>Assume</u>

- Documents have a unique id and a path
- A document is a sequence of words

- Document d = $[w_1, w_2, ..., w_n]$

- Search for a single, exact-match word
 - Does document D have word w?
 - The list of documents **D** that have word w?



Advanced Issues

- More expressive query semantics
 - Multi-word
 - Locality: ["what is it"] vs. ["what it is"] vs. ["what", "is", "it"]
- Ranked results
 - Document-ranking algorithm (e.g. PageRank)
 - Efficient ranked retrieval
- Dynamics
 - Document addition/removal/modification
 - Rank
 - Document changes
 - Integration of real-time variables (e.g. location)



Modeling Semantic Memory

- Semantic memory is a human's long-term store of facts about the world, independent of the context in which they were originally learned
- The ACT-R (<u>http://act-r.psy.cmu.edu</u>) model of semantic memory has been successful at explaining a variety of psychological phenomena (e.g. retrieval bias, forgetting)
- The model does *not* scale to large memory sizes, which hampers complex experiments



1 December 2014

Scale Fail [AFRL '09]

Retrieval Latency: Chunks in DM x Retrieval Constraints x Type of DM (Error Bars: 95% Confidence Interval)



Inverted Index

Memory Representation



- Document = Node
- Word = edge

Example cue: last(obama), spouse(X)



Ranking [Anderson et al. '04]





Example

Semantic Objects: Features





1 December 2014

Derbinsky

Fall 2014

Inverted Index



Index Statistics



Top-1 Non-Ranked Retrieval





Derbinsky

Introducing Rank



Wentworth Institute of Technology

COMP570 – Database Applications

Derbinsky

Fall 2014

Ranked Retrieval Algorithm #1 Sort on Query

Inverted Index



Wentworth Institute of Technology

COMP570 – Database Applications

Derbinsky

Fall 2014

Ranked Retrieval Algorithm #3 Static Sort

Inverted Index



24

Wentworth Institute of Technology

COMP570 – Database Applications

Derbinsky

Fall 2014

Ranked Retrieval Algorithm #2 Static Sort

Inverted Index



Hybrid Approach

- Empirically supported cardinality threshold, θ
- If (cardinality > θ): Sort on Query [#1]
 - Candidate enumeration scales with # of objects with large cardinality (empirically rare)
- If (cardinality $\leq \theta$): Static Sort [#2]
 - Bias updates must be locally efficient
 - Objects affected: O(1)
 - Computation: O(1)



Some Results

Inverted index (via SQLite) + new approach was **fast** and **scaled**!

>30x faster than off-the-shelf database (on >3x data)!



Takeaways

- A common approach to large-scale search is indexing: using data structure(s) to improve access speed
- An inverted index is commonly used for full-text search (even in situations that might not look like it)
- Inverted indexes are fast, scalable, and straightforward to implement
- Know your indexes/data structures! Careful problem analysis and algorithm development can often beat generic approaches
 - Even if you don't use a DBMS, DBMS methods can be very useful in a variety of applications!

