Information Retrieval WS 2012 / 2013

Lecture 3, Wednesday November 7th, 2012 (List intersection: fancy algorithms + lower bounds)

> Prof. Dr. Hannah Bast Chair of Algorithms and Data Structures Department of Computer Science University of Freiburg

Overview of this lecture

Organizational

- Your experiences with Exercise Sheet 2 (Ranking)
- List intersection
 - Key algorithm in every search engine
 - Some algorithm engineering tips
 - Fancier algorithms + runtime analysis
 - Matching lower bound
 - Exercise Sheet 3:

Implement the asymptotically optimal algorithm Is it really faster than the simple linear-time algorithm? Some optional theoretical tasks, good for exam preparation BURG

Experiences with ES2 (ranking)

Summary / excerpts last checked November 7, 15:43

- Feasible in reasonable time for most
- Manual relevance assessment is quite time-consuming
- Wiki table for comparing results please, like last semester
 Ok, we will have one again for this exercise !
- Some of you are quite passionate about the coding
 Great, I hope you become addicted to it !
- Lecture pen too thick / too many colors ... let's see today
- How to write unit tests ... I'll live code an example today
- More coding, less theory please ... ok, ok
- Great recordings / streamed from iPhone to TV ... yeah!

Your results for ES2 (ranking)

Some interesting observations

- Binary (k = 0, b = 0) was clearly the worst

– In the tf.idf list, some linguists and other theory guys

Many occurrences of either "relativity" or "theory" can boost the score

 Albert Einstein was top for tf.idf, but not even in the top-10 for BM25 with standard settings

BM25 favors short articles that are mainly about rel. theo.

To get AE at the top, a global document score (like the PageRank for web pages) would be more appropriate ... later lecture

Let's first revisit the linear-time algorithm ...

- ... and see what we can improve there
- Java: ArrayList<Integer> much worse than native int[]
- C++: vector<int> is as good as int[] with option -O3
- Minimize the number of branches within small-bodied loops
- Note 1: hard to predict / understand some of the effects
 To understand what is really going on, look at the machine code (in C++) or byte code (Java) ... see seminar Java vs. C++ on Wiki
- Note 2: when comparing algorithms, always repeat each run at least 3 times, to make caching effects transparent

BURG

List Intersection 2/5

Algorithmic improvement 1

- Call the smaller list A, and the longer list B
- Search the elements from A in B, using binary search

both souted as usual

k = #elements in A, n = #elements in B

- This has time complexity $\Theta(k \cdot \log n)$







List Intersection 5/5

Time complexity of this algorithm

– Let d_1 , ..., d_k be the gaps between the locations of the k elements of A in B

 d_1 = from beginning to first location

- Note that $\Sigma_i d_i \leq n$ = the number of elements in B
- Then the time complexity is $O(\Sigma_i \log d_i)$
- Goal: find a formula that is independent of the d_i
- Idea: maximize $\Sigma_i \log d_i$ under the constraint $\Sigma_i d_i \leq n$
- This is called optimization with side constraints or Lagrangian optimization \dots next slide \mathcal{RESULT} \mathcal{L} \mathcal{L}



• A heuristic approach to increase performance

- Idea: place pointers at "strategic" positions in the list, which allow to "skip" large parts during intersection
- Question: where to place how many pointers?



Let's first look at both union and intersection

- For list **union**, the simple linear-time algorithm has a time complexity of $\Theta(k + n)$

This is optimal, because the result of the union contains k + n elements, and every algorithm has to output them

- For list **intersection**, the exponential-binary-search algorithm has a time complexity of $\Theta(k \cdot \log (n/k))$

Exercise (optional): prove that $k \cdot \log (n/k)$ never worse than n + k (and usually much better, e.g. for k = const)

But maybe we can do even better?

(Note: the output size is at most k for intersection)

BURG

Recall: lower bound for comparison-based sorting

- There are n! possible outputs
- The algorithm has to distinguish between all of them
- Each comparison distinguishes between two cases
- Hence we need at least $\log_2 n!$ comparisons
- By Stirling's formula $(n/e)^n \le n! \le n^n$
- Hence $\log_2 n! \sim n \cdot \log n$
- Hence every comparison-based sorting algorithm has a running time of $\Omega(n \cdot \log n)$
- Note: not true for non-comparison based algorithms
 For example, 0-1 sequences can be sorted by counting

JNI FREIBURG Similar argument for intersecting two lists A and B

- As before, k = #elements in A, n = #elements in B

- How many different ways are there to locate the k elements from A within the n elements from B
- Observation: each such way corresponds to a tuple $(j_1, ..., j_k)$ where $0 \le j_1 \le ... \le j_k \le n$

 j_i is simply the location of the i-th element of A in B location 0 means before the first element location i > 0 means after the i-th element

How many such tuples are there? THE ANSWER: $\binom{M+k}{2} \ge \binom{M}{2}^{k}$ $\log\left(\frac{M}{2}\right)^{k} = 2 \cdot \log\frac{M}{2}$

There is a similar quantity which is easy to count

- The number of tuples $(i_1, ..., i_k)$ where $1 \le i_1 < ... < i_k \le m$
- This is just the number of size-k subsets of $\{1, ..., m\}$
- Their number is just m over $k = m! / (k! \cdot (m k)!) \ge (m/k)^k$
- Let us relate this to the number we are interested in: The number of tuples $(i_1, ..., i_k)$ where $0 \le i_1 \le ... \le i_k \le n$

Now it's easy to derive the lower bound

- By the previous arguments, the number of ways to locate the k elements from A in the n element of B is m over $k \ge (m/k)^k$, where m = n + k
- With the same argument as in the sorting lower bound, any comparison-based algorithm hence has to make $\log_2 (m \text{ over } k) \ge \log_2 (m/k)^k$ comparisons
- This is $k \cdot \log_2(1+n/k)$ and hence $\Omega(k \cdot \log_2(n/k))$

17

REIBURG

References

In the Raghavan/Manning/Schütze textbook

Section 2.3: Faster intersection with skip pointers

Relevant Papers

A simple algorithm for merging two linearly ordered setsF.K. Hwang and S. LinSICOMP 1(1):31-39, 1980

A fast set intersection algorithm for sorted sequences R. Baeza-Yates CPM, LNCS 3109, 31–39, 2004

Relevant Wikipedia articles

http://en.wikipedia.org/wiki/Lagrange multiplier

