Lecture 9: Linear Sorting

Steven Skiena

Department of Computer Science State University of New York Stony Brook, NY 11794–4400

http://www.cs.stonybrook.edu/~skiena

Topic: Problem of the Day

Problem of the Day

The *nuts and bolts* problem is defined as follows. You are given a collection of n bolts of different widths, and n corresponding nuts. You can test whether a given nut and bolt together, from which you learn whether the nut is too large, too small, or an exact match for the bolt. The differences in size between pairs of nuts or bolts can be too small to see by eye, so you cannot rely on comparing the sizes of two nuts or two bolts directly. You are to match each bolt to each nut.

- 1. Give an $O(n^2)$ algorithm to solve the nuts and bolts problem.
- 2. Suppose that instead of matching all of the nuts and bolts, you wish to find the smallest bolt and its corresponding nut. Show that this can be done in only 2n 2 comparisons.
- 3. Match the nuts and bolts in expected $O(n \log n)$ time.



Topic: Lower Bounds on Sorting

Can we sort in $o(n \lg n)$ **?**

Any comparison-based sorting program can be thought of as defining a decision tree of possible executions. Running the same program twice on the same permutation causes it to do exactly the same thing, but running it on different permutations of the same data causes a different sequence of comparisons to be made on each.



Claim: the height of this decision tree is the worst-case complexity of sorting.

Lower Bound Analysis

Since any two different permutations of n elements requires a different sequence of steps to sort, there must be at least n!different paths from the root to leaves in the decision tree. Thus there must be at least n! different leaves in this binary tree.

Since a binary tree of height h has at most 2^h leaves, we know $n! \le 2^h$, or $h \ge \lg(n!)$. By inspection $n! > (n/2)^{n/2}$, since the last n/2 terms of the

product are each greater than n/2. Thus

 $\log(n!) > \log((n/2)^{n/2}) = n/2\log(n/2) \rightarrow \Theta(n\log n)$

Stirling's Approximation

By Stirling's approximation, a better bound is $n! > (n/e)^n$ where e = 2.718.

$$h \ge \lg(n/e)^n = n \lg n - n \lg e = \Omega(n \lg n)$$



Topic: Bucketsort: Non-Comparison-Based Sorting

Non-Comparison-Based Sorting

All the sorting algorithms we have seen assume binary comparisons as the basic primative, questions of the form "is x before y?".

But how would you sort a deck of playing cards?

Most likely you would set up 13 piles and put all cards with the same number in one pile.

With only a constant number of cards left in each pile, you can use insertion sort to order by suite and concatenate everything together.

If we could find the correct pile for each card in constant time, and each pile gets O(1) cards, this algorithm takes O(n) time.

Bucketsort

Suppose we are sorting n numbers from 1 to m, where we know the numbers are approximately uniformly distributed. We can set up n buckets, each responsible for an interval of m/n numbers from 1 to m



Given an input number x, it belongs in bucket number $\lceil xn/m \rceil$.

If we use an array of buckets, each item gets mapped to the right bucket in O(1) time.

Bucketsort Analysis

With uniformly distributed keys, the expected number of items per bucket is 1. Thus sorting each bucket takes O(1) time!

The total effort of bucketing, sorting buckets, and concatenating the sorted buckets together is O(n).

What happened to our $\Omega(n \lg n)$ lower bound!

Worst-Case vs. Assumed-Case

Bad things happen to bucketsort when we assume the wrong distribution.



We might spend linear time distributing our items into buckets and learn *nothing*.

Problems like this are why we worry about the worst-case performance of algorithms!

Real World Distributions

The worst case "shouldn't" happen if we understand the distribution of our data.

Consider the distribution of names in a telephone book.

- Will there be a lot of Skiena's?
- Will there be a lot of Smith's?
- Will there be a lot of Shifflett's?

Either make *sure* you understand your data, or use a good worst-case or randomized algorithm!

The Shifflett's of Charlottesville

For comparison, note that there are seven Shifflett's (of various spellings) in the 1000 page Manhattan telephone directory.

Shifflett Rebble & Duckerwille	092.7957	Chillion Island	2219 Williamenter Mil
SURFICE ACTAL & VANALITIE	703-1731	Semiller Caller	
Shifflett Debra S SR 617 Quinque	985-8813	Shimett James	2 801 Stonenenge AV
Chiffield Dalma CD400	9025-360	Chiffinit James	C Stanardenille
Sumer Admie Subha	703-3000	Stutter Annuas	A Briter Astrine
Shifflett Delmax Crozet	823-5901	Shifflett James	E Earlysville
Chieflate Plannan G. Manihuma		Childland Innen	E .Is EES Clausiand Av
Suuner nember a wenitm		Summer Same	E OL SOT MELINING MA
100 Graanbrier Ter	973-7195	Shifflett James	F & Lois LongMeadow
Childheld Denies Di (07 D.L.	000-0003	Chiffield James	E & Mannall Did 71
SUMMENT VERISE KT 62/ Dyke	782-907/	SUMMENT OBUILDE	L C ACLINCK MINAT
Shifflett Dennis Stanardiville	985-456A	Shifflett James	J 1430 Rugby Av
CLIMICAL DATA IS DURING CONTINUE		Childland Isman	W Da Casses An
SUILLET DEUDE IN STRUCTURE	A93.5AX4	Shiller Asines	V OF GROUDE WA
Shifflett Dewey F 21667	025-6576	Shifflett James	L SR33 Stanardsville
	AND - 78/0	Children In Income	
SWITHERT DEWEY U DYRE	7037/297	· ?Willingt 4410422	Carlystant
Chifflett Diana 508 Rainbridge Av	070-70tc	Shifflett James	O Stanardsville •••••
Alithian Cabu C Baadala and			B Ald Lunghburg Bd
SUMMENT DODY & MAILICIA KID	280-4227	Printiger County	IK ON LYNCHOURS HE **
Chilliett Doub(Ola Dt 671	074.7442	Chifflatt James	R 81753 Exmont
	77777799	Maleistass admenute	

Non-Comparison Sorts Don't Beat the Bound

Radix sort works by partitioning on the lowest order characters first, maintaining this order to break ties. It takes time O(nm) to sort n strings of length m, or time *linear* in the input size! But m must be $\Omega(\log n)$ before the strings are all distinct! Sorting n arbitrary, distinct keys cannot be done better than $\Theta(n \log n)$.

