CS61B Lecture #25		Shell's sort
Today: Sorting, cont.		Idea: Improve insertion sort by first sorting <i>distant</i> elements: • First sort subsequences of elements $2^k - 1$ apart:
• Standard methods		- sort items $\#(-2^k - 1) - 2(2^k - 1) - 3(2^k - 1)$ then
 Froperties of standard methods Selection 		- sort items #2, $2 + 2^k - 1$, $2 + 2(2^k - 1)$, $3(2^k - 1)$,, then - sort items #2, $2 + 2^k - 1$, $2 + 2(2^k - 1)$, $2 + 3(2^k - 1)$,, then
Readings for Today: DS(IJ), Chapter 8;		- etc. - sort items # $2^k - 2$, $2(2^k - 1) - 1$, $3(2^k - 1) - 1$,,
Readings for Next Topic: Balanced searches,	DS(IJ), Chapter 9;	– Each time an item moves, can reduce #inversions by as much as $2^k + 1$.
		• Now sort subsequences of elements $2^{k-1}-1$ apart:
		- sort items #0, $2^{k-1} - 1$, $2(2^{k-1} - 1)$, $3(2^{k-1} - 1)$,, then - sort items #1, $1 + 2^{k-1} - 1$, $1 + 2(2^{k-1} - 1)$, $1 + 3(2^{k-1} - 1)$,, -:
		• End at plain insertion sort ($2^0 = 1$ apart), but with most inversions gone.
		$ullet$ Sort is $\Theta(N^{1.5})$ (take CS170 for why!).
Last modified: Thu Oct 28 16:21:30 2004	C561B: Lecture #25 1	Last modified: Thu Oct 28 16:21:30 2004 C561B: Lecture #25 2
Example of Shell's Sort		Sorting by Selection: Heapsort
	#I #C	Idea: Keep selecting smallest (or largest) element.
15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0	120 1	 Really bad idea on a simple list or vector.
		 But we've already seen it in action: use heap.
		• Gives $O(N \lg N)$ algorithm (N remove-first operations).
0 14 13 12 11 10 9 8 7 6 5 4 3 2 1 15	91 10	 Since we remove items from end of heap, we can use that area to accumulate result:
0 7 6 5 4 3 2 1 14 13 12 11 10 9 8 15	42 20	original: 19 0 -1 7 23 2 42 heapified: 42 23 19 7 0 2 -1
	4 19	23 7 19 -1 0 2 42 19 7 2 -1 0 23 42 7 0 2 -1 19 23 42
<u>0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15</u> T: Tryersions left	U -	
C: Comparisons needed to sort subsequences.		0 -1 2 7 19 23 42 -1 0 2 7 19 23 42

Merge Sorting

Idea: Divide data in 2 equal parts; recursively sort halves; merge results.

- Already seen analysis: $\Theta(N \lg N)$.
- Good for *external* sorting:
 - First break data into small enough chunks to fit in memory and sort.
 - Then repeatedly merge into bigger and bigger sequences.
 - Can merge K sequences of arbitrary size on secondary storage using $\Theta(K)$ storage.
- For internal sorting, can use *binomial comb* to orchestrate:

Illustration of Internal Merge Sort

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)



Last modified: Thu Oct 28 16:21:30 2004

CS61B: Lecture #25 6

Quicksort: Speed through Probability

Idea:

Last modified: Thu Oct 28 16:21:30 2004

- Partition data into pieces: everything > a pivot value at the high end of the sequence to be sorted, and everything \leq on the low end.
- Repeat recursively on the high and low pieces.
- For speed, stop when pieces are "small enough" and do insertion sort on the whole thing.
- Reason: insertion sort has low constant factors. By design, no item will move out of its will move out of its piece [why?], so when pieces are small, #inversions is, too.
- Have to choose pivot well. E.g.: *median* of first, last and middle items of sequence.

Example of Quicksort

- In this example, we continue until pieces are size ≤ 4 .
- Pivots for next step are starred. Arrange to move pivot to dividing line each time.
- Last step is insertion sort.

1	.6 10	13 1	8 -4	-7 12	-5	19 1	5 0	22	29 3	4 -1*	
-	4 -5	-7	-1	8 13	12 10) 19	15	0 22	29	34 16) *
-	4 -5	-7	-1	5 13 1	.2* 10	0	16	19*	22 2	9 34	18
-	4 -5	-7	-1 1	0 0	12	15 1	3 10	5 18	19	29	34 22
• Now ev	/eryt	hing	is "cl	ose ta	o" rig	ht, s	so ju	st d	o ins	sertic	on sort:

-7 -5 -4 -1 0 10 12 13 15 16 18 19 22 29 34

CS61B: Lecture #25 5

Performance of Quicksort

• Probabalistic time:

Last modified: Thu Oct 28 16:21:30 2004

- If choice of pivots good, divide data in two each time: $\Theta(N\lg N)$ with a good constant factor relative to merge or heap sort.
- If choice of pivots bad, most items on one side each time: $\Theta(N^2).$
- $\Omega(N \lg N)$ in best case, so insertion sort better for nearly ordered input sets.
- \bullet Interesting point: randomly shuffling the data before sorting makes $\Omega(N^2)$ time very unlikely!

13 31 21 -4 37 4* 11 10 39 2 0 40 59 51 49 46 60

-4 0 2 4 37 13 11 10 39 21 31* 40 59 51 49 46 60

-4 0 2 4 21 13 11 10 31 39 37 40 59 51 49 46 60

-4 0 2 4 21 13 11 10 31 37 39 40 59 51 49 46 60

Looking for #6 to right of pivot 4:

Looking for #1 to right of pivot 31:

Just two elements; just sort and return #1:

Quick Selection

The Selection Problem: for given k, find $k^{\dagger h}$ smallest element in data.

- Obvious method: sort, select element #k, time $\Theta(N \lg N)$.
- \bullet If $k \leq$ some constant, can easily do in $\Theta(N)$ time:
 - Go through array, keep smallest \boldsymbol{k} items.
- \bullet Get $\textit{probably}\,\Theta(N)$ time for all k by adapting quicksort:
 - Partition around some pivot, $p, \, {\rm as}$ in quicksort, arrange that pivot ends up at dividing line.
 - Suppose that in the result, pivot is at index m, all elements \leq pivot have indicies $\leq m.$
 - If m = k, you're done: p is answer.
 - If m > k, recursively select k^{th} from left half of sequence.
 - If m < k , recursively select $(m-k-1)^{\mbox{th}}$ from right half of sequence.

Selection ExampleSelection PerformanceProblem: Find just item #10 in the sorted version of array:• For this algorithm, if m roughly in middle each time, cost isInitial contents:
 $51 160 21 - 4 37 4 49 10 40^{+} 59 0 13 2 39 11 46 31$
0
Looking for #10 to left of pivot 40:• For this algorithm, if m roughly in middle each time, cost is $C(N) = \begin{cases} 1, & \text{if } N = 1, \\ N + C(N/2), & \text{otherwise.} \\ = N + N/2 + \ldots + 1 \\ = 2N - 1 \in \Theta(N) \end{cases}$

- \bullet But in worst case, get $\Theta(N^2),$ as for quicksort.
- \bullet By another, non-obvious algorithm, can get $\Theta(N)$ worst-case time for all k (take CS170).

Result: 39

CS61B: Lecture #25 9

Last modified: Thu Oct 28 16:21:30 2004

CS61B: Lecture #25 10

Better than N lg N?		Beyond Comparison: Di	stribution Counting			
• Can prove that if all you can do to keys is compar	e them then sorting	• But suppose can do more than comp	oare keys?			
 must take Ω(N lg N). Basic idea: there are N! possible ways the in 	nput data could be	• For example, how can we sort a set range from 0 to kN , for some smal	t of N integer keys whose values I constant k ?			
scrambled.		• One technique: <i>count</i> the number of	of items $< 1, < 2$, etc.			
 Therefore, your program must be prepared to do binations of move operations. 	$\sim N!$ different com-	• If $M_p = \#$ items with value $< p$, the with value $= p$ must be $\#M = i$	en in sorted order, the j^{th} item			
 Therefore, there must be N! possible combinations of outcomes of all the if tests in your program (we're assuming that comparisons are 2-way). 		with value p must be $\#M_p + j$. • Gives linear-time algorithm.				
 Since each if test goes two ways, number of pos- comes for k if tests is 2^k. 	sible different out-					
• Thus, need enough tests so that $2^k > N!$, which it	means $k \in \Omega(\lg N!)$.					
 Using Stirling's approximation, 						
$m! \in \sqrt{2\pi m} \left(\frac{m}{e}\right)^m \left(1 + \Theta\left(\frac{1}{m}\right)\right)$,					
this tells us that						
$k \in \Omega(N \lg N).$						
Last modified: Thu Oct 28 16:21:30 2004	C561B: Lecture #25 13	Last modified: Thu Oct 28 16:21:30 2004	CS61B: Lecture #25 14			
Distribution Counting Examp	le	Radix S	ort			
• Suppose all items are between 0 and 9 as in this	example:	Idea: Sort keys one character at a t	time.			
7 0 4 0 9 1 9 1 9 5 3 7 3 1 4		 Can use distribution counting for each digit. 				
		• Can work either right to left (LSD	radix sort) or left to right (MSD			
3 3 1 2 2 1 1 3 0 3	Counts	radix sort)				
		• LSD radix sort is venerable: used t	for punched cards.			
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Running sum	Initial: set, cat, cad, con	ı, bat, can, be, let, bet			
0 0 0 1 1 1 2 3 3 4 4 5 6 7 7 0 3 6 9 11 12 13	7 7 9 9 9 16	bet let bat (by char #2) be cad con set 'i' 'd' 'n' 't' be cad con can set cat bat let bet	bat bet Pass 2 cat let (by char #1) cad be con 'a' 'e' 'o' cad can cat bat be set let bet con			
• "Counts" line gives # occurrences of each key.						
$ullet$ "Running sum" gives cumulative count of keys \leq e	each value		con bet cat			
 which tells us where to put each key: 		Pass 3 (by char #0)	be can bat cad let set			
• The first instance of key k goes into slot m, whe of key instances that are < k.	ere m is the number	bat, be	'b' 'c' 'l' 's' e, bet, cad, can, cat, con, let, set			

CS61B: Lecture #25 16

MSD Radix Sort

• A bit more complicated: must keep lists from each step separate

Α

* set, cat, cad, con, bat, can, be, let, bet

* bat, be, bet / cat, cad, con, can / let / set

bat / \star be, bet / cat, cad, con, can / let / set

bat / be / bet / * cat, cad, con, can / let / set

bat / be / bet / * cat. cad. can / con / let / set 2

bat / be / bet / cad / can / cat / con / let / set

• But, can stop processing 1-element lists

Performance of Radix Sort

- Radix sort takes $\Theta(B)$ time where B is total size of the key data.
- Have measured other sorts as function of #records.
- How to compare?
- \bullet To have N different records, must have keys at least $\Theta(\lg N)$ long [why?]
- Furthermore, comparison actually takes time $\Theta(K)$ where K is size of key in worst case [why?]
- So $N \lg N$ comparisons really means $N(\lg N)^2$ operations.
- While radix sort takes $B = N \lg N$ time.
- On the other hand, must work to get good constant factors with radix sort.

Last modified: Thu Oct 28 16:21:30 2004	CS61B: Lecture #25 17	Last modified: Thu Oct 28 16:21:30 2004	CS61B: Lecture #25 18

And Don't Forget Search Trees

Idea: A search tree is in sorted order, when read in inorder.

- Need balance to really use for sorting [next topic].
- \bullet Given balance, same performance as heapsort: N insertions in time $\lg N$ each, plus $\Theta(N)$ to traverse, gives

$$\Theta(N + N \lg N) = \Theta(N \lg N)$$

Summary

 \bullet Insertion sort: $\Theta(Nk)$ comparisons and moves, where k is maximum amount data is displaced from final position.

- Good for small datasets or almost ordered data sets.

- Quicksort: $\Theta(N\lg N)$ with good constant factor if data is not pathological. Worst case $O(N^2).$
- \bullet Merge sort: $\Theta(N\lg N)$ guaranteed. Good for external sorting.
- \bullet Heapsort, treesort with guaranteed balance: $\Theta(N\lg N)$ guaranteed.
- \bullet Radix sort, distribution sort: $\Theta(B)$ (number of bytes). Also good for external sorting.

posn

0

1

2

1