

## 算法设计与分析 Algorithms Design & Analysis

### 第四讲：线性时间的排序法

1

## 线性时间的排序法(Sorting in Linear Time)

- Lower bound for comparison-based sorting(基于比较的排序方法的开销下限)
- Counting sort(记数排序)
- Radix sort(基数排序法)
- Bucket sort(桶排序)

## 排序法回顾(Sorting So Far)

- Insertion sort(插入排序):
  - Easy to code(易于实现)
  - Fast on small inputs (less than ~50 elements)(适合小规模输入)
  - Fast on nearly-sorted inputs(对近似排序好的输入效果好)
  - $O(n^2)$  worst case(最坏情况)

3

## 排序法回顾(Sorting So Far)

- Merge sort(合并排序):
  - Divide-and-conquer(分治):
    - Split array in half(分割)
    - Recursively sort subarrays(递归处理)
    - Linear-time merge step(线性时间合并)
  - $O(n \lg n)$  worst case(最坏情况)

4

## How Fast Can We Sort?(我们能否排得更快?)

- Insertion sort  $\Theta(n^2)$
- Merge sort  $\Theta(n \lg n)$
- What is common to all these algorithms? (上述算法共同点?)
  - These algorithms sort by making comparisons between the input elements (基于元素之间的两两比较)

5

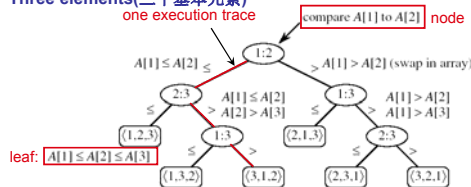
## 比较排序(Comparison Sorting)

- Comparison sorts use comparisons between elements to gain information about an input sequence  $\langle a_1, a_2, \dots, a_n \rangle$  (通过比较获得输入元素之间的相对顺序, 比较的操作有: )
- $a_i < a_j, a_i \leq a_j, a_i = a_j, a_i \geq a_j, \text{ or } a_i > a_j$
- To sort  $n$  elements, comparison sorts must make  $\Omega(n \lg n)$  comparisons in the worst case (在最坏情况, 基于比较思想排序的必须要做  $\Omega(n \lg n)$  次比较)

6

## 判定树模型(Decision Tree Model)

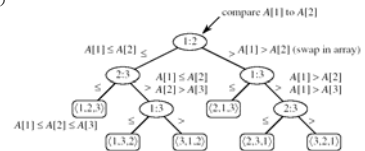
- Represents the comparisons made by a sorting algorithm on an input of a given size: models all possible execution traces(表示基于比较排序算法的所有可能比较情况)
- Control, data movement, other operations are ignored(忽略了其他操作,比如控制,数据转移等操作)
- Count only the comparisons(只考虑比较)
- Three elements(三个基本元素)



7

## Decision Tree Model

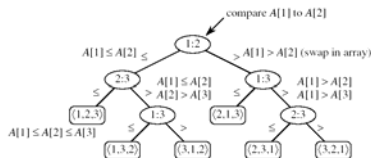
- All permutations on  $n$  elements must appear as one of the leaves in the decision tree  $n!$  permutations
- $n!$ 种排列对应同数目的叶节点
- Worst-case number of comparisons (最坏情况)
  - the length of the longest path from the root to a leaf (从根到叶最长的路径)



8

## 判定树模型(Decision Tree Model)

- Goal: finding a lower bound on the running time on any comparison sort algorithm(目的是为了确定基于比较的排序算法的下限)
  - find a lower bound on the heights of all decision trees in which each permutation appears as a reachable leaf (判定树中, 所有元素的排列情况作为叶节点出现, 目的是确定判定树的高度)



9

## 引理 (Lemma)

- Any binary tree of height  $h$  has at most  $2^h$  leaves (任何树高为 $h$ 的二分树至多有 $2^h$ 个叶节点)

**Proof:** induction on  $h$  (证明, 归纳法)

**Basis:**  $h = 0 \Rightarrow$  tree has one node, which is a leaf (只有叶节点)  
 $2^0 = 1$

**Inductive step:** assume true for  $h-1$  (假设对于 $h-1$ 成立)

- Extend the height of the tree with one more level
- Each leaf becomes parent to two new leaves (当高度增加1时, 即 $h-1 \rightarrow h$ , 原来的叶节点都成为2个新叶节点的父节点)

$$\begin{aligned} \text{No. of leaves at level } h &= 2 * (\text{no. of leaves at level } h-1) \\ &= 2 * 2^{h-1} \\ &= 2^h \end{aligned}$$

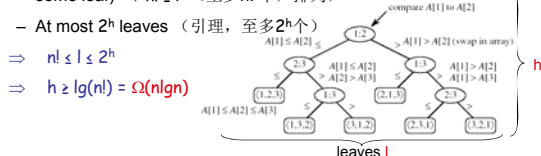
10

## 比较排序下限 (Lower Bound for Comparison Sorts)

**Theorem:** Any comparison sort algorithm requires  $\Omega(n \lg n)$  comparisons in the worst case. (任何基于比较的排序算法在最坏情况下至少需要做 $\Omega(n \lg n)$ 次比较)

**Proof:** How many leaves does the tree have? (叶节点的数目)

- At least  $n!$  (each of the  $n!$  permutations if the input appears as some leaf)  $\Rightarrow n! \leq I$  (至少 $n!$ 个, 排列)
- At most  $2^h$  leaves (引理, 至多 $2^h$ 个)



We can beat the  $\Omega(n \lg n)$  running time if we use other operations than comparisons! (放弃比较的排序思想, 可以获得更优开销, 突破 $\Omega(n \lg n)$ , 但是是有条件的)

11

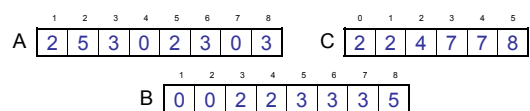
## 计数排序 (Counting Sort)

**Assumption:**

- The elements to be sorted are integers in the range 0 to  $k$  (将要排序的元素是0到 $k$ 之间的整数)

**Idea:** (算法思想)

- Determine for each input element  $x$ , the number of elements smaller than  $x$  (对于每一个元素 $x$ , 判断小于或等于元素 $x$ 的元素的个数)
- Place element  $x$  into its correct position in the output array (然后根据计数的情况将元素放入到相应的位置)



12

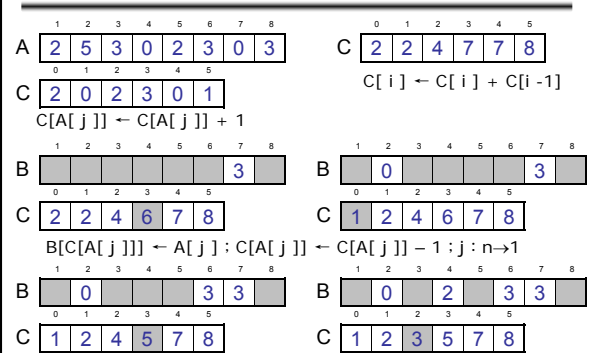
## 计数排序 (COUNTING-SORT)

**Alg.:** COUNTING-SORT(A, B, n, k)

1. **for**  $i \leftarrow 0$  **to**  $k$
2.     **do**  $C[i] \leftarrow 0$
3.     **for**  $j \leftarrow 1$  **to**  $n$
4.         **do**  $C[A[j]] \leftarrow C[A[j]] + 1$
5.      $\triangleright C[i]$  contains the number of elements equal to  $i$  ( $C[i]$  包含等于  $i$  的元素的个数)
6.     **for**  $i \leftarrow 1$  **to**  $k$
7.         **do**  $C[i] \leftarrow C[i] + C[i-1]$
8.      $\triangleright C[i]$  contains the number of elements  $\leq i$  ( $C[i]$  包含小于或等于  $i$  的元素的个数)
9.     **for**  $j \leftarrow n$  **downto**  $1$
10.         **do**  $B[C[A[j]]] \leftarrow A[j]$
11.          $C[A[j]] \leftarrow C[A[j]] - 1$

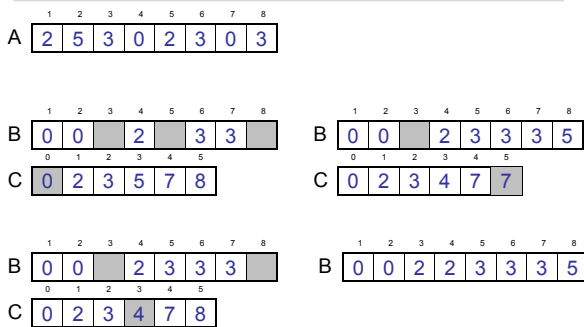
13

## Example



14

## Example (cont.)



15

## 算法分析 (Analysis of Counting Sort)

**Alg.:** COUNTING-SORT(A, B, n, k)

1. **for**  $i \leftarrow 0$  **to**  $k$   $\left. \begin{array}{l} 1. \\ 2. \end{array} \right\} \Theta(k)$
2.     **do**  $C[i] \leftarrow 0$
3.     **for**  $j \leftarrow 1$  **to**  $n$   $\left. \begin{array}{l} 3. \\ 4. \end{array} \right\} \Theta(n)$
4.         **do**  $C[A[j]] \leftarrow C[A[j]] + 1$
5.      $\triangleright C[i]$  contains the number of elements equal to  $i$
6.     **for**  $i \leftarrow 1$  **to**  $k$   $\left. \begin{array}{l} 6. \\ 7. \end{array} \right\} \Theta(k)$
7.         **do**  $C[i] \leftarrow C[i] + C[i-1]$
8.      $\triangleright C[i]$  contains the number of elements  $\leq i$
9.     **for**  $j \leftarrow n$  **downto**  $1$   $\left. \begin{array}{l} 9. \\ 10. \\ 11. \end{array} \right\} \Theta(n)$
10.         **do**  $B[C[A[j]]] \leftarrow A[j]$
11.          $C[A[j]] \leftarrow C[A[j]] - 1$

Overall time:  $\Theta(n + k)$

16

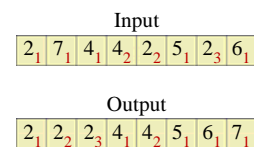
## 分析 (Analysis of Counting Sort)

- Overall time:  $\Theta(n + k)$
- In practice we use COUNTING sort when  $k = O(n)$   
 $\Rightarrow$  running time is  $\Theta(n)$  (小范围整数的排序)
- Counting sort is **stable** (稳定的)
  - Numbers with the same value appear in the same order in the output array (相等的输入元素的位置先后顺序在排列后输出中的先后顺序不变)

17

## Stable的定义

- A sorting algorithms is **stable** if for any two indices  $i$  and  $j$  with  $i < j$  and  $a_i = a_j$ , element  $a_i$  precedes element  $a_j$  in the output sequence. (稳定的定义)



## 基数排序(Radix Sort)

- Extending the Range of Sortable Integers
- (扩大了可排序的整数范围)

## 基数排序思想(The Idea of Radix Sort)

- Any number can be represented in radix-k notation.(数的基表示)
- Examples:**
  - $16536 = 1 \cdot 10^4 + 6 \cdot 10^3 + 5 \cdot 10^2 + 3 \cdot 10^1 + 6 \cdot 10^0$
  - $\langle 101101 \rangle_2 = 1 \cdot 2^5 + 0 \cdot 2^4 + 1 \cdot 2^3 + 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 45$

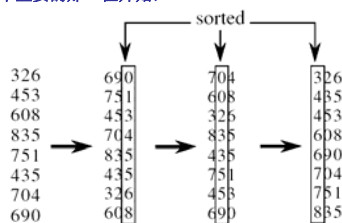
## 基数排序 (RADIX-SORT)

Alg.: RADIX-SORT(A, d)

for  $i \leftarrow 1$  to  $d$

do use a stable sort to sort array A on digit i (用稳定的排序方法对第i位排序)

- 1 is the lowest order digit, d is the highest-order digit (从最后, 也就是最不重要的那一位开始)



21

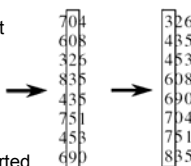
## Analysis of Radix Sort

- Given  $n$  numbers of  $d$  digits each, where each digit may take up to  $k$  possible values, RADIX-SORT correctly sorts the numbers in  $\Theta(d(n+k))$  ( $n$ 个 $d$ 位的数, 以 $k$ 为基)
  - One pass of sorting per digit takes  $\Theta(n+k)$  assuming that we use counting sort (如果用计数排序, 排一位需要 $\Theta(n+k)$ )
  - There are  $d$  passes (for each digit) (排 $d$ 位需要 $\Theta(d(n+k))$ )

22

## 基数排序的正确性 (Correctness of Radix sort)

- We use induction on number of passes through each digit (归纳法)
- Basis:** If  $d = 1$ , there's only one digit, trivial (1位, 显然)
- Inductive step:** assume digits  $1, 2, \dots, d-1$  are sorted (1, 2, ..., d-1已排序)
  - Now sort on the  $d$ -th digit (三种情况)
    - If  $a_d < b_d$ , sort will put  $a$  before  $b$ : correct  
 $a < b$  regardless of the low-order digits
    - If  $a_d > b_d$ , sort will put  $a$  after  $b$ : correct  
 $a > b$  regardless of the low-order digits
    - If  $a_d = b_d$ , sort will leave  $a$  and  $b$  in the same order and  $a$  and  $b$  are already sorted on the low-order  $d-1$  digits (等于, 低位已经排好) WHY?



23

## 桶排序 (Bucket Sort)

- Assumption:**
  - the input is generated by a random process that distributes elements uniformly over  $[0, 1)$  (输入是在 $[0, 1)$ 上一致均匀分布)
- Idea:**
  - Divide  $[0, 1)$  into  $n$  equal-sized buckets (将 $[0, 1)$   $n$ 等分, 称为空桶)
  - Distribute the  $n$  input values into the buckets (将 $n$ 个输入分派到这些空桶之中)
  - Sort each bucket (每个桶内排序)
  - Go through the buckets in order, listing elements in each one (按顺序遍历每个桶中的元素)
- Input:**  $A[1 \dots n]$ , where  $0 \leq A[i] < 1$  for all  $i$
- Output:** elements  $a_i$  sorted
- Auxiliary array:**  $B[0 \dots n-1]$  of linked lists, each list initially empty (桶, 初始为空)

24

## 桶排序算法 (BUCKET-SORT)

*Alg.*: BUCKET-SORT( $A, n$ )

for  $i \leftarrow 1$  to  $n$

do insert  $A[i]$  into list  $B[\lfloor nA[i] \rfloor]$  (注意下标)

for  $i \leftarrow 0$  to  $n - 1$

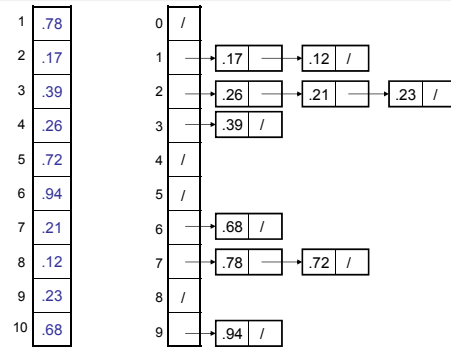
do sort list  $B[i]$  with insertion sort (桶内插入排序)

concatenate lists  $B[0], B[1], \dots, B[n-1]$   
together in order

return the concatenated lists

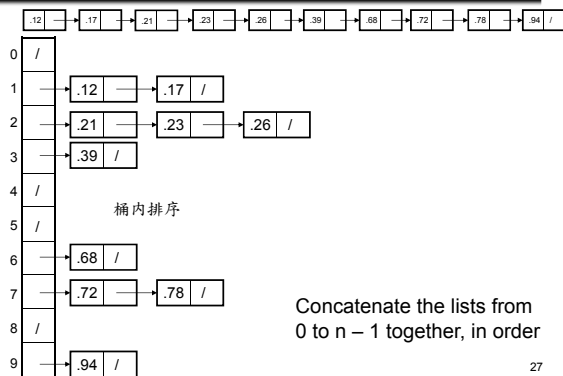
25

## Example - Bucket Sort



26

## Example - Bucket Sort



27

## 桶排序的正确性(Correctness of Bucket Sort)

- Consider two elements  $A[i], A[j]$  (考虑两元素  $A[i], A[j]$ )
- Assume without loss of generality that  $A[i] \leq A[j]$  (不失一般性可假设  $A[i] \leq A[j]$ )
- Then  $\lfloor nA[i] \rfloor \leq \lfloor nA[j] \rfloor$  (有  $\lfloor nA[i] \rfloor \leq \lfloor nA[j] \rfloor$ )
  - $A[i]$  belongs to the same group as  $A[j]$  or to a group with a lower index than that of  $A[j]$  ( $A[i]$  要么与  $A[j]$  在一个桶中, 要么在比  $A[j]$  更低的桶中)
- If  $A[i], A[j] \leq \lfloor nA[j] \rfloor$  belong to the same bucket: (同一桶中, 桶内排序保障  $A[i] \leq A[j]$  顺序正常)
  - insertion sort puts them in the proper order
- If  $A[i], A[j]$  are put in different buckets: (不同桶中, 连接过程保障正确)
  - concatenation of the lists puts them in the proper order

28

## 平均情况(Average-Case Running Time of Bucket Sort)

- Lemma:** Given that the input sequence is drawn uniformly at random from  $[0, 1)$ , the average-case running time of Bucket Sort is  $O(n)$ . (输入满足一致性随机分布, 桶排序的平均运行开销为  $O(n)$ . 证明可阅读教材)

## 结论(Conclusion)

- Any comparison sort will take at least  $n \lg n$  to sort an array of  $n$  numbers (基于比较:  $n \lg n$ )
- We can achieve a better running time for sorting if we can make certain assumptions on the input data: (对输入作更细致的考虑)
  - Counting sort: each of the  $n$  input elements is an integer in the range 0 to  $k$  (计数排序: 0到 $k$ 的整数)
  - Radix sort: the elements in the input are integers represented with  $d$  digits (基数排序:  $d$ 整数位)
  - Bucket sort: the numbers in the input are uniformly distributed over the interval  $[0, 1)$  (桶排序:  $[0, 1)$ 一致性分布)

30