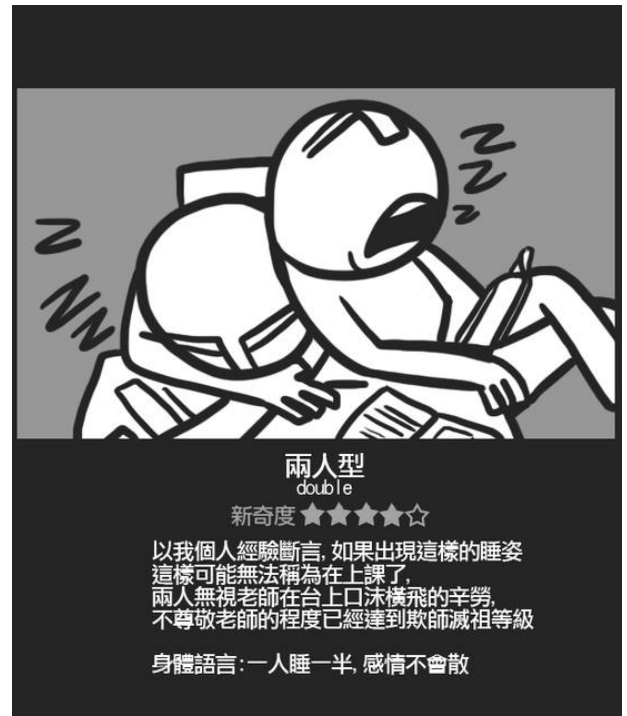


SORTING IN LINEAR TIME

Michael Tsai

2012/5/22



Sorting Algorithm in Linear Time?

- 上回講過, Comparison-based (比較元素然後交換其順序) sorting的complexity lower bound為 $\Omega(n \log n)$
- 如何突破這個障礙?
- 不要使用comparison-based的方法
- 因此可達到 $O(n)$
- 需要利用額外的假設
- 通常是“以空間換取時間”

例子：電話排序問題

- 電話排序假設全部電話號碼有 s 個種可能 (所有排列組合)
- 我們有 n 組電話要排
- 則使用超大表格排序法需要 $O(s)$ 的時間 (空間換取時間)

Input: {20000002, 89999999, 20000000, ...}

共 s 個

Index	Array
20000000	有
空	空
20000002	有
...	...
...	...
89999999	有

例子：電話排序問題

- 如果使用comparison-based sorting平均需要 $O(n \log n)$
- 假設超大表格排序法所需時間為 $c_1 s$
- 假設comparison-based排序法所需時間為 $c_2 n \log n$
- 假設 $10c_1 = c_2$
- break even point為 $s = 10 n \log n$
- 假設以台北市電話號碼為例, $s = 10^8$
- $10^8 = 10 n \log n, n < 10^8$
- break even point大約為 $n = k \times 10^6$
- 當 n 大於此數則non-comparison sorting比較快

Counting Sort

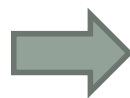
假設: 知道可能出現的所有input種類
個數 K (而且 $K = O(n)$)

Input array

	1	2	3	4	5	6	7	8
A	2	5	3	0	2	3	0	3

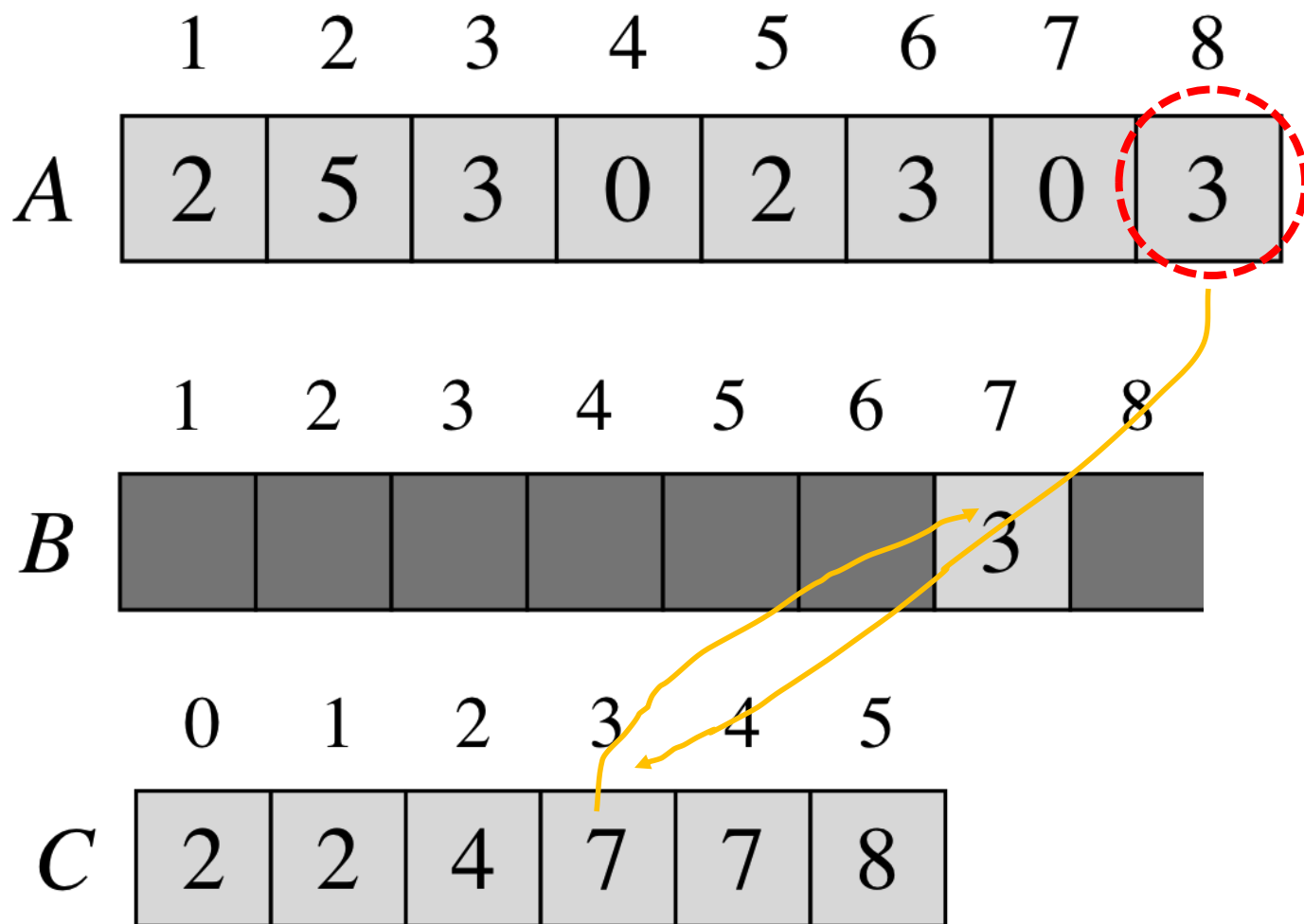
	0	1	2	3	4	5
C	2	0	2	3	0	1
	0	1	2	3	4	5
C	2	2	4	7	7	8

Counting: 算每個數字出現了幾次



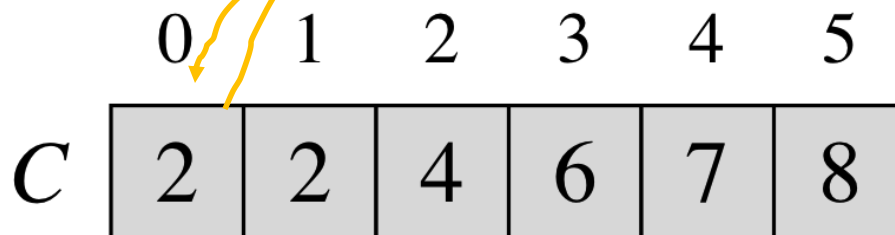
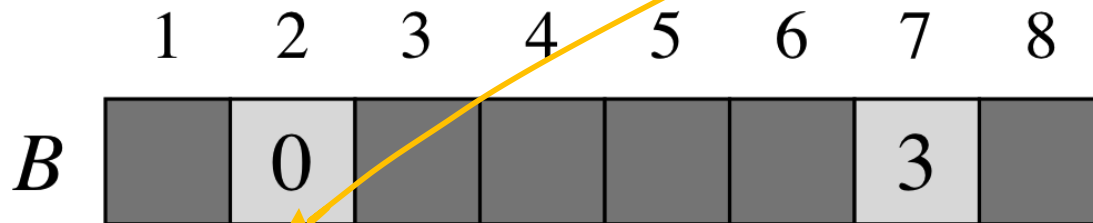
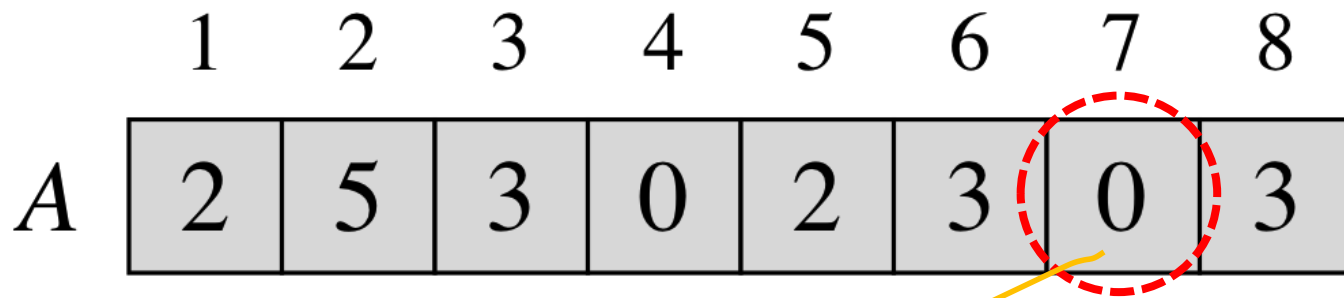
“Cumulative Mass Function”
算小於或等於index的數字出現了幾次

Counting Sort



比3小或相等的有7個

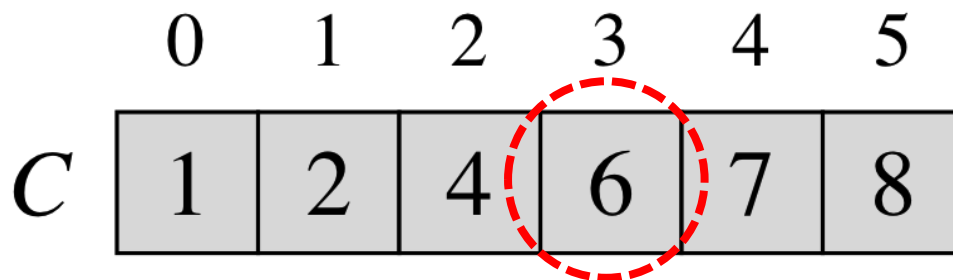
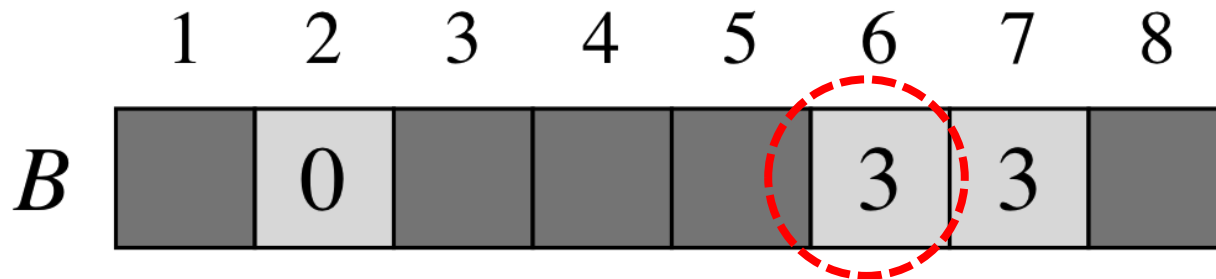
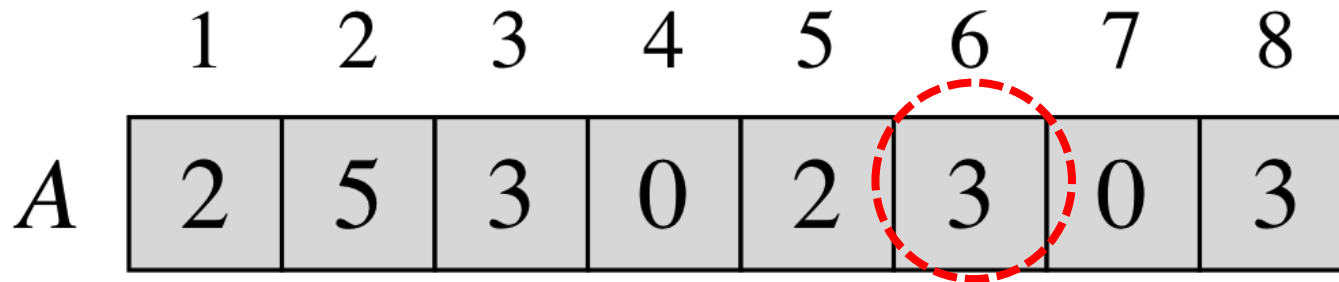
Counting Sort



比0小或相等的有2個

注意Counting Sort為stable sort. 因為是從最後面依序放入output array, 因此同樣大小的元素會依原本input array中順序放入.

Counting Sort



Counting Sort

A: input array
 n: input總個數
 B: output array
 K: 可能出現的input element總數

```
void CountingSort (int A[], int n, int B[], int K) {
    int C[K], i, j;
    for (i=0; i<K; i++)
        C[i]=0
    for (j=0; j<n; j++)
        C[A[j]] = C[A[j]] + 1;
    for (i=1; i<K; i++)
        C[i] = C[i] + C[i-1];
    for (j=n-1; j>=0; j--) {
        B[C[A[j]]] = A[j];
        C[A[j]]--;
    }
}
```

$O(K)$

$O(n)$

$O(K)$

$O(n)$

Counting: 數每種element共幾個

"CMF": 數比每種element小或相等的共幾個

從input array最後一個開始依序放入output array

$O(n + K) = O(n)$

Sorting on several keys

- 假設 K 有很多個sub-key

- $K = (K_1, K_2, \dots, K_r)$

Most significant key \uparrow Least significant key \uparrow

- 則 $K_x \leq K_y$ iff
- $K_{x,i} = K_{y,i}, 1 \leq i \leq j$ and $K_{x,j+1} < K_{y,j+1}$ for some $j < r$, or
- $K_{x,i} = K_{y,i}, 1 \leq i \leq r$
- 則我們可以有以下的sorting 方法.
- Most Significant Digit first (MSD) sorting
 - 先依照most significant key sort, 然後依序往least significant key sort過去
- Least Significant Digit first (LSD) sorting
 - 先依照least significant key sort, 然後依序往most significant key sort過去

Radix Sort

假設:有“多個key”

以個位數排序

329	720
457	355
657	436
839	457
436	657
720	329
355	839

以十位數排序

720
329
436
839
355
457
657

注意:
十進位一樣的部分
會照前一回合的順序
(也就是個位數的順序!)

因此用來排序的演算法
必須是stable!

以百位數排序

329
355
436
457
657
720
839

先依照most significant digit sort, 然後依序往least significant key sort
過去: Most Significant Digit first (MSD) sorting

Radix Sort

如果先依照least significant key sort, 然後依序往most significant key sort過去: Least Significant Digit first (LSD) sorting?

以百位數排序

329	329
457	355
657	457
839	436
436	657
720	720
355	839

以十位數排序

329
720
436
839
355
457
657

正確的做法必須切分成小的subset再去排序! (因此LSD比較麻煩)

329
355
436
457
657
720
839

如果直接使用接下來的十位數排序會亂掉!

Radix Sort

以個位數排序

329
457
657
839
436
720
355

以十位數排序

720
355
436
457
657
329
839

以百位數排序

329
355
436
457
657
720
839

$O(n)$

$O(n)$

$O(n)$

假設使用Counting Sort在每一回合做排序

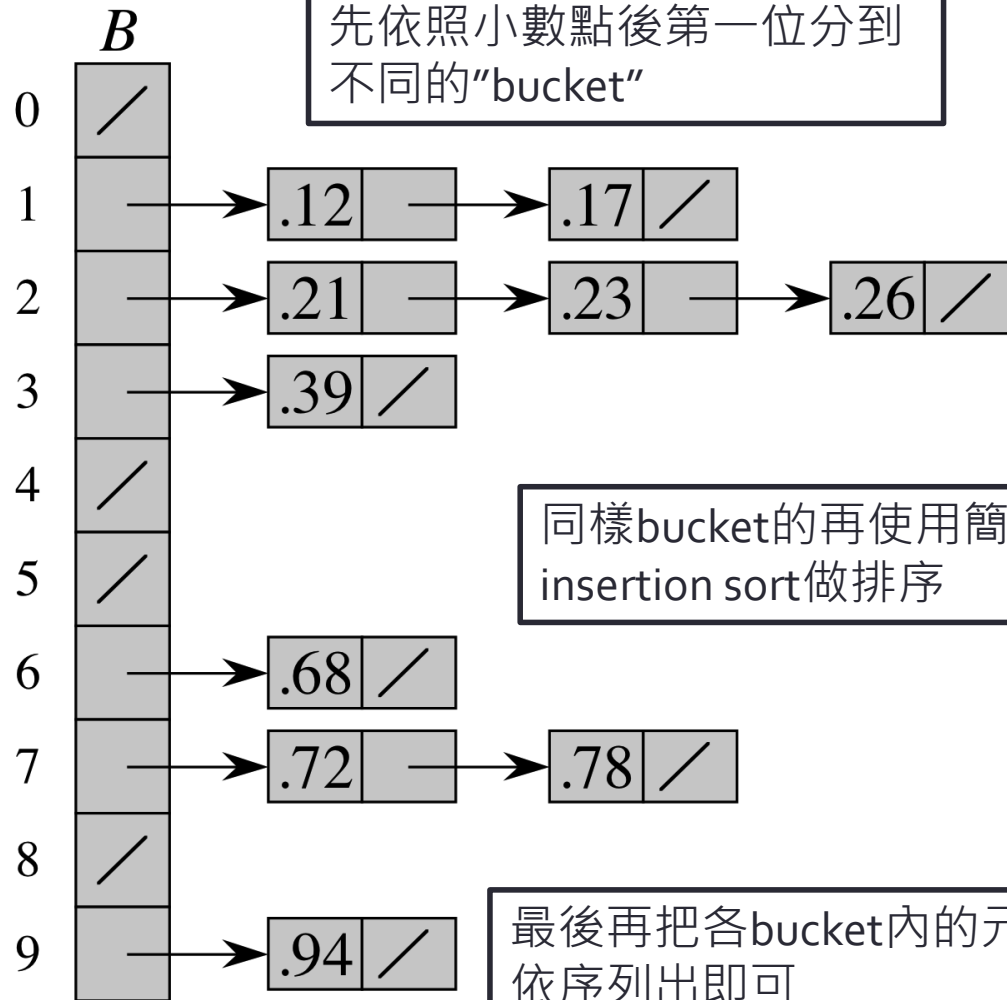
d個, d=位數

$O(dn)$

Bucket Sort

假設: input是從uniform distribution取出來的

	<i>A</i>
1	.78
2	.17
3	.39
4	.26
5	.72
6	.94
7	.21
8	.12
9	.23
10	.68



Bucket Sort

- 假設: input是從uniform distribution取出來的
- 因此分到各個bucket的數量會差不多.

- 大略計算一下所需的時間:

走過每個bucket的時間

- $T(n) = \Theta(n) + \sum_{i=0}^{n-1} O(n_i^2)$

總時間

每個bucket內insertion sort所需時間

- $$E[T(n)] = E\left[\Theta(n) + \sum_{i=0}^{n-1} O(n_i^2)\right]$$

$$= \Theta(n) + \sum_{i=0}^{n-1} O(E[n_i^2]) = \Theta(n) + n \cdot O\left(2 - \frac{1}{n}\right) = \Theta(n)$$

$$E[n_i^2] = 2 - \frac{1}{n}$$

(證明過程請見Cormen p202-203)

Today's Reading Assignment

- Cormen 8.2, 8.3, 8.4