

Σε κάθε αλγόριθμο οι ασυμπτωτικές εκτιμήσεις χρόνου είναι πάντοτε βασισμένες στο μοντέλο υπολογισμού RAM (random access machine). Στο μοντέλο υπολογισμού RAM η προσπέλαση της μνήμης γίνεται μέσω υπολογισμού διευθύνσεων, δηλαδή υποστηρίζεται η χρήση πινάκων (αττάς).

Σε αντιδιαστολή με τη RAM υπάρχει και η μηχανή δείκτην (pointer machine) όπου η προσπέλαση μνήμης είναι δυνατή μόνο μέσω δείκτην, δηλαδή οι διευθύνσεις είναι τιμές δείκτην. Η pointer machine εκπροσωπεί υλοποιήσεις αλγορίθμων με records και είναι ασθενέστερο μοντέλο απ' τη RAM.

ΔΙΑΤΑΞΗ ΣΤΟΙΧΕΙΩΝ (SORTING)

- 1) Bubblesort: κεντρική ιδέα: α) αλλάξε τα ζεύγη εκείνα που βρίσκονται εκτός διατάξης
β) επανέλαβε το προηγούμενο βήμα έως ότου διαταχθούν όλα τα ζεύγη.

Εμφανίζεται σε δύο μορφές. Στην πρώτη μορφή η επεξεργασία των στοιχείων λαμβάνει χώρα από αριστερά προς τα δεξιά και στη δεύτερη μορφή η επεξεργασία των στοιχείων γίνεται από αριστερά προς τα δεξιά και απ' τα δεξιά προς τα αριστερά εναλλάξ (cocktail-shakesort).

π.χ. input: 4 3 8 6 2 5 7

```

      ↓
3 4 8 6 2 5 7
      ↓
3 4 6 8 2 5 7
      ↓
3 4 6 2 8 5 7
      ↓
3 4 6 2 5 8 7
      ↓
3 4 6 2 5 7 8
      ↓
      ⋮
  
```

output: 2 3 4 5 6 7 8

Αρχικά, συγκρίνουμε το 1^ο στοιχείο με το 2^ο στοιχείο. Αν το 2^ο είναι μικρότερο του 1^{ου} τότε τα ανταλλάσσουμε και μετά συγκρίνουμε το πυχόν νέο 2^ο στοιχείο με το 3^ο κι αν $2^o > 3^o$ τα αλλάζουμε, αλλιώς προχωράμε στο 4^ο στοιχείο το οποίο συγκρίνουμε με το 5^ο, κ.ο.κ. Μόλις θα έχουμε διατρέξει 1 φορά το input αττάς, το μεγαλύτερο στοιχείο θα βρίσκεται στα δεξιότερα και το πρόβλημα θα έχει βελτιωθεί κατά 1. Μετά καθούμε να εφαρμόσουμε την ίδια διαδικασία με την προηγούμενη για τα υπόλοιπα $n-1$ στοιχεία (αν είναι η αρχικός αριθμός των στοιχείων του input).

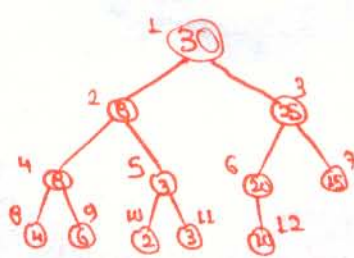
Πλουσιμότητα: $n + (n-1) + (n-2) + \dots + 1 = \frac{n(n+1)}{2} = O(n^2)$ στη χειρότερη περίπτωση, γιατί έχουμε n βήματα με αριθμό συγκρίσεων $\leq n-1$ στο καθένα.

2) Heapsort:

- Μια δομή heap (σωρός) είναι ένα δένδρο του οποίου οι κόμβοι φέρουν αριθμούς με την ιδιότητα ότι για κάθε κόμβο v ισχύει: αριθμός(ρ(v)) ≥ αριθμός(v). Σε ένα τέτοιο δένδρο ο κάθε πατέρας είναι μεγαλύτερος απ' τα παιδιά του και ο μεγαλύτερος αριθμός βρίσκεται στη ρίζα.
- Στην υλοποίηση του heapsort χρησιμοποιούνται ισοζυγισμένα και δυϊκά heaps, όπου δυϊκό σημαίνει ότι κάθε κόμβος έχει ούχ 2 παιδιά με ενδεχόμενη εξαίρεση κάποιων κόνων και ισοζυγισμένα σημαίνει ότι τα φύλλα έχουν βάθος k ή $k+1$ ή ότι βρίσκονται στο k ή $k+1$ επίπεδο. Τα φύλλα βάθους $k+1$ είναι συμπληρωμένα προς τα αριστερά.
- Το ισοζυγισμένο και δυϊκό heap μπορεί ονομαστικά να παρασταθεί με ένα αττάς.

Π.Χ.

30	8	25	8	3	20	15	4	6	2	3	10
1	2	3	4	5	6	7	8	9	10	11	12



Η αρίθμηση των κόμβων γίνεται απ' τη ρίζα και οι:

$$\text{father}(i) = \left\lfloor \frac{i}{2} \right\rfloor$$

$$\text{left}(i) = 2i$$

$$\text{right}(i) = 2i + 1$$

(μόνο για ισοδυστομένο και διατεταγμένο heap)

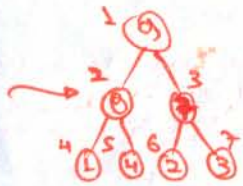
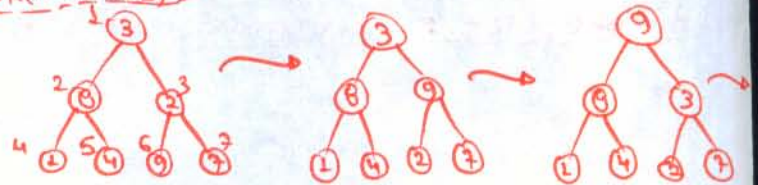
Ο heap sort χωρίζεται σε δύο φάσεις. Τη φάση δόμησης και τη φάση διαλογής.

Στη φάση δόμησης διαμορφώνεται το αντίστροφο heap ενώ στην φάση διαλογής χρησιμοποιείται η $\text{extractmax}()$, η οποία όταν εφαρμόζεται σε ένα κόμβο βρίσκει το μεγαλύτερο παιδί του και το βάζει σαν πατέρα, ενώ ο πατέρας βγαίνει απ' το heap. (διαλέγει και απομακρύνει το μεγαλύτερο στοιχείο, μετατρέπει τα υπόλοιπα στο αντίστροφο heap και επαναλαμβάνει τη διαδικασία μέχρι τέλους)

Π.Χ.

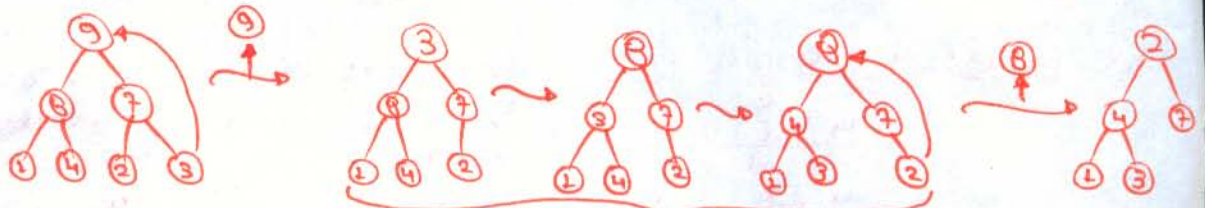
3	8	2	1	4	9	7
1	2	3	4	5	6	7

φάση δόμησης:

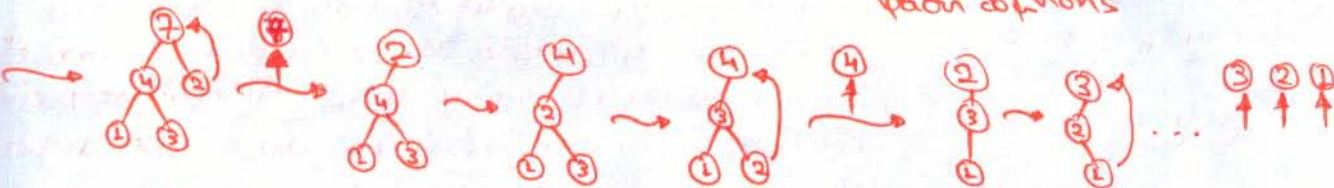


Κάθε φορά κάνουμε εναλλαγή των αριθμών για τους οποίους δεν ικανοποιείται η heap-property.

φάση διαλογής:



φάση δόμησης



Στη φάση διαλογής παίρνουμε το μέγιστο στοιχείο (το οποίο είναι αποθηκευμένο στη ρίζα) και στη θέση του βάζουμε το πιο δεξί απ' τα παιδιά. Στη συνέχεια χρησιμοποιείται η extractmax η οποία εξάγει το μεγαλύτερο στοιχείο. Έπειτα εφαρμόζεται η φάση δόμησης ώστε να ικανοποιείται η heap-property και επαναλαμβάνεται πάλι η ίδια διαδικασία. Τελικά, όλα τα στοιχεία που θα εξαχθούν απ' το αρχικό heap θα είναι διατεταγμένα σε ένα πίνακα.

Πλομολογώμενα για τη φάση δόμησης: $\sum_{i=1}^n (\text{κόστος πρόσδεσης του } i[i]) = \sum_{i=1}^n O(h(i)) = O(n + \sum_{k=0}^{\infty} k \cdot (\text{πλήθος κόμβων } i \text{ με } h(i)=k)) = O(n + \sum_{k=0}^{\infty} k \cdot \frac{n}{2^k}) = O(n \cdot (1 + \sum_{k=0}^{\infty} \frac{1}{2^k})) = O(n)$ *το πλήθος των κόμβων i με $h(i)=k$ είναι $\leq \frac{n}{2^k}$

Πλομολογώμενα για τη φάση διαλογής: το μέγιστο κάθε φορά στοιχείο απομακρύνεται απ' τη ρίζα του heap το πολύ η φορές. Μετά την απομάκρυνση κάθε μέγιστου, η μετατροπή σε heap κοστίζει $O(\log n)$. Οπότε το συνολικό κόστος είναι $O(n \log n)$. (στη χειρότερη περίπτωση)

Ο heap sort είναι ο καλύτερος αλγόριθμος για διατάξη η στοιχείων και έχει πλομολογώμενα $O(n \log n)$, είναι δηλαδή βέλτιστος (optimal). Ισχύει: $\Pi(X)$ (σπουδαιότερη αλγ. ταξ.) \geq ύψος δένδρου απόφασης

• Αν για η στοιχεία φτιάξουμε ένα δένδρο απόφασης, τότε ο αριθμός των εισόδων είναι $n!$ και άρα θα έχουμε $n!$ φύλλα. Ξέραμε, όμως, ότι ένα δυϊκό δένδρο ύψους h έχει 2^h φύλλα. Οπότε, για ένα διατεταγμένο αλγόριθμο που χρησιμοποιεί δυϊκό δένδρο, θα ισχύει: $2^h \geq n! \Rightarrow h \geq \log_2 n! \Rightarrow h \geq \log_2 n^n \Rightarrow h \geq n \log_2 n$ ($\log_2 n! = \sum_{i=1}^n \log_2 i \geq \int_1^n \log_2 x dx = \Omega(n \log n)$)

3) QuickSort

Βασίζεται στην αρχή divide-and-conquer: Για ένα πίνακα $A[p...r]$ εφαρμόζονται τα εξής:

Divide: Ο πίνακας $A[p...r]$ διαχωρίζεται σε δύο μη αδειούς υποπίνακες $A[p...q]$ και $A[q+1...r]$ τέτοιους ώστε κάθε στοιχείο του $A[p...q]$ να είναι μικρότερο ή ίσο από κάθε στοιχείο του $A[q+1...r]$. Το στοιχείο q υπολογίζεται ως μέρος της συνάρτησης διαχωρισμού.

Conquer: Οι δύο υποπίνακες $A[p...q]$ και $A[q+1...r]$ διατάσσονται με αναδρομικές κλήσεις του quicksort.

Combine: Αφού ο κάθε πίνακας είναι διατεταγμένος τότε συνενώνονται τους θα μας δώσει τον ταξινομημένο πίνακα $A[p...r]$.

Ο αλγόριθμος είναι ο εξής:

```

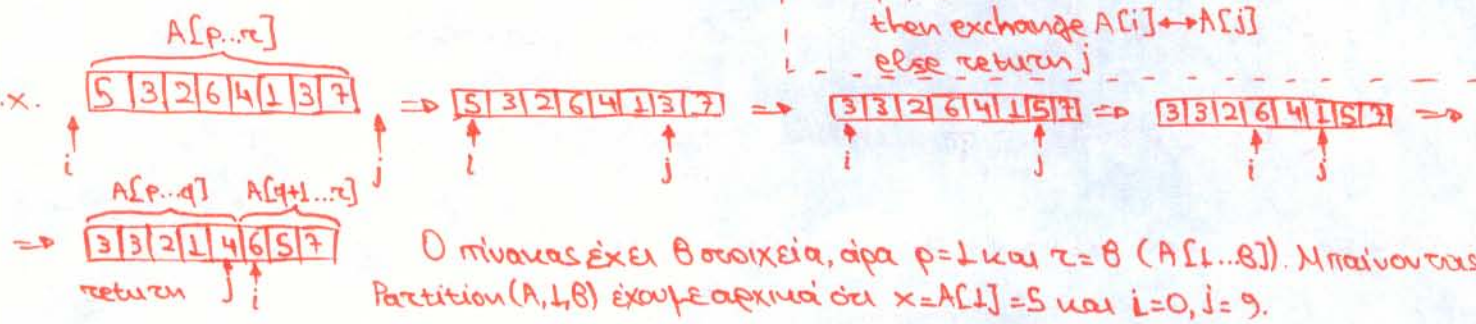
QuickSort(A, p, r)
if p < r
    then q ← Partition(A, p, r)
         QuickSort(A, p, q)
         QuickSort(A, q+1, r)
    
```

```

Partition(A, p, r)
x ← A[p]
i ← p-1
j ← r+1
while TRUE
    do repeat j ← j-1
      until A[j] ≤ x
    repeat i ← i+1
      until A[i] ≥ x
    if i < j
      then exchange A[i] ↔ A[j]
    else return j
    
```

3

Για ολόκληρο πίνακα η κλήση είναι QuickSort(A, 1, length[A])



~~Όταν το $j=7$ τότε το $A[j]=A[7]=3 \leq 5$ και τότε το $i=1$, αφού $A[i]=5 \geq 5$. Οπότε τα 3 και 5 ανταλλάσσονται, αφού $i < j$ ($1 < 7$). Η ίδια διαδικασία επαναλαμβάνεται έως ότου $j < i$ οπότε και επιστρέψουμε, καλώντας αναδρομικά τον αλγόριθμο για τα κομμάτια του πίνακα (υποπίνακες) $A[1...5]$ και $A[6...8]$.~~

Όταν το $j=7$ τότε το $A[j]=A[7]=3 \leq 5$ και τότε το $i=1$, αφού $A[i]=5 \geq 5$. Οπότε τα 3 και 5 ανταλλάσσονται, αφού $i < j$ ($1 < 7$). Η ίδια διαδικασία επαναλαμβάνεται έως ότου $j < i$ οπότε και επιστρέψουμε, καλώντας αναδρομικά τον αλγόριθμο για τα κομμάτια του πίνακα (υποπίνακες) $A[1...5]$ και $A[6...8]$.

Πομπητικότητα χειρότερης περίπτωσης: Αν ορίσουμε σαν $Q(n)$ το μέγιστο πλήθος των συγκρίσεων που εκτελεί ο QuickSort σε έναν πίνακα μήκους n τότε ισχύει η παρακάτω ανισότητα:

$$Q(n) \leq n+1 + \max_{1 \leq j \leq n} \{Q(j-1) + Q(n-j)\}, \text{ όπου ισχύει } Q(0) = Q(1) = 0 \text{ και επαγωγικά}$$

για τον 1ο διαχωρισμό σε 2 υποπληθύνσεις

αποδεικνύεται ότι $Q(n) \leq \frac{(n+1)(n+2)}{2} - 3 = O(n^2)$

Πομπητικότητα μέσης περίπτωσης:

Έχουμε τις υποθέσεις:
 • Όλα τα στοιχεία του πίνακα $A[1...n]$ είναι ανά δύο διαφορετικά
 • Έχουμε ομοιόμορμη κατανομή, δηλαδή καθένα των $n!$ δυνατών μεταθέσεων των δεδομένων είναι ισοπίθανο.

Εάν για τον πίνακα S έχουμε ότι το $S[i]=k$ τότε τα υποπρόβλήματα μεγέθους $k-1$ και $n-k$ επιλύονται πάλι τις προηγούμενες υποθέσεις. Για το υποπρόβλημα μεγέθους $k-1$ η απόδειξη είναι προφανής.

• Με τον 1^ο διαμερισμό του quicksort παίρνουμε το σχηματισμό: $K = k_1 \dots k_{n-1} k_n \dots k_n$ όπου $k_1 \dots k_{n-1}$ είναι μετάθεση των $1, \dots, n-1$ στοιχείων του πίνακα και $k_n \dots k_n$ είναι μετάθεση των $k+1, \dots, n$ αντίστοιχα. Αν είναι τ το πλήθος των εναλλαγών στην πραγματοποιηθείσα σειρά, δηλαδή μεναξύ των αριθμών $1, \dots, n-1$ και $k+1, \dots, n$, τότε εναλλάσσονται τ ζευγάρια. Οπότε, υπάρχουν:

$\sum_{\tau \geq 0} \binom{n-1}{\tau} \cdot \binom{n-k}{\tau}$ ακολουθίες που μπορούν να υπάρχουν στην ακολουθία K μετά το διαμερισμό.

Θεωρώντας, τώρα, ότι το μέσο πλήθος των συγκρίσεων που εμπεριέχει ο Quicksort σε ένα πίνακα μήκους n είναι $\overline{QS}(n)$, θα ισχύει $\overline{QS}(0) = \overline{QS}(1) = 0$, αφού αυτές τις χρησιμοποιούμε υποθέσεις συμπεραίνουμε ότι $\text{prob}(S[1]=k) = 1/n$.

Επομένως, θα έχουμε:

$$\overline{QS}(n) = \underbrace{(n+1)}_{\text{για το διαμερισμό}} + \frac{1}{n} \cdot \sum_{s=1}^n [\overline{QS}(s-1) + \overline{QS}(n-s)] \Rightarrow \overline{QS}(n) = (n+1) + \frac{2}{n} \sum_{s=0}^{n-1} \overline{QS}(s)$$

$$\left\{ \begin{array}{l} n \cdot \overline{QS}(n) = n \cdot (n+1) + 2 \sum_{s=0}^{n-1} \overline{QS}(s) \\ (n+1) \cdot \overline{QS}(n+1) = (n+1)(n+2) + 2 \sum_{s=0}^n \overline{QS}(s) \end{array} \right\} \xrightarrow{(-)} (n+1)\overline{QS}(n+1) - n\overline{QS}(n) = (n+1)(n+2) - n(n+1) + 2\overline{QS}(n) \Rightarrow$$

(για $s=n$)

$$\Rightarrow (n+1)\overline{QS}(n+1) = (n+2)\overline{QS}(n) + 2(n+1) \Rightarrow \overline{QS}(n+1) = 2 + \frac{n+2}{n+1} \cdot \overline{QS}(n) = 2 + \frac{n+2}{n+1} \left[2 + \frac{n+1}{n+1} \overline{QS}(n+1) \right]$$

$$= 2 \cdot \left[1 + (n+2) \cdot \left(\frac{1}{n+1} + \frac{1}{n} + \dots + \frac{1}{2} \right) \right] \Rightarrow \overline{QS}(n) = 2 \cdot \left[1 + (n+1) \cdot \left(\frac{1}{n} + \frac{1}{n-1} + \dots + \frac{1}{2} \right) \right] \Rightarrow$$

$$\Rightarrow \overline{QS}(n) = 2 + 2(n+1) \cdot \sum_{i=1}^n \frac{1}{i} - 2(n+1) = 2(n+1) \cdot \sum_{i=1}^n \frac{1}{i} - 2n$$

Όμως το $\sum_{i=1}^n \frac{1}{i}$ είναι αρμονικός αριθμός και ισχύει $\sum_{i=1}^n \frac{1}{i} \leq \int_1^n \frac{1}{x} dx = 1 + \ln(n)$

Οπότε: $\overline{QS}(n) = 2(n+1) + 2(n+1)\ln(n) - 2n \Rightarrow \overline{QS}(n) = 2(n+1)\ln(n) + 2 \Rightarrow \overline{QS}(n) = O(n \log n)$

Άρα ο Quicksort διατάσσει n στοιχεία σε χειρότερο χρόνο $O(n^2)$ και μέσο χρόνο $O(n \log n)$.

⊗ Δίνεται ένα σύνολο αριθμών από 0 έως n , όπου το $n = 2^m - 1$, οι οποίοι μπορούν να αναπαρασταθούν με m bits. Οι αριθμοί αναπαριστούνται στο δυαδικό σύστημα αριθμών. Ένας απ' αυτούς λείπει. Να βρεθεί ο αριθμός σε γραμμικό χρόνο.

• Κάθε αριθμός του συνόλου αναπαρίσταται με m bits: $x = \underbrace{x_1}_{MSB} x_2 \dots \underbrace{x_m}_{LSB}$, όπου το κάθε ψηφίο $x_i \in \{0, 1\}$. Ο χρόνος που απαιτείται για να προσπελασθούν τα ψηφία είναι $O(m)$ και για όλους τους αριθμούς είναι $O(n \cdot m)$, ενώ αν έχουμε λογαριθμικό μέτρο θα είναι $O(n \cdot \log n)$.

Μπορούμε να εφαρμόσουμε μια πιο έξυπνη λύση που βασίζεται στην αρχή divide and conquer. Αν π.χ. θεωρήσουμε ότι $m=3$ τότε το σύνολο των αριθμών είναι $n = 2^3 - 1 = 7$, άρα από 0...7. Απ' αυτούς λείπει ένας, τον οποίο και θα βρούμε:

1	0	0
0	0	0
0	0	1
0	1	0
1	1	1
1	0	1
0	1	1

1	0	0
0	0	0
0	1	0
0	1	0

⊙ 10

Ξεκινάμε με τα LSB's και συγκρίνουμε τους άσσους με τα μηδενικά. Βλέπουμε ότι λείπει ένα μηδενικό. Άρα στην επόμενη φάση παίρνουμε τους αριθμούς με 0 στο τέλος, ενώ ένα 0 θα προστεθεί στο τελευταίο bit του αριθμού που λείπει. Στη 2^η φάση βλέπουμε ότι λείπει ένας 1, οπότε συνεχίζοντας με τον ίδιο τρόπο έχουμε τον αριθμό 110.

Επιμένω, η πολυπλοκότητα θα είναι $O(n)$ στο 1^ο βήμα, $O(\frac{n}{2})$ στο 2^ο, $O(\frac{n}{4})$ στο 3^ο ... $O(1)$ στο τελευταίο ($2^m - 1$ φορές)

External sorting

- Θαυρούμε ότι διαθέτουμε ένα file μήκους N , μια κίβρα μνήμη μέγεθους M και $2p$ βοηθητικές ταινίες (external devices). Τελικά, στην έξοδο θα έχουμε το file διατεταγμένο.
- Αρχικά, διαχωρίζουμε τις βοηθητικές ταινίες σε p -input και p -output ταινίες
- Στην κίβρα μνήμη παράγονται διατεταγμένα υποfiles μήκους M .
- Αναμειγνύουμε τα διατεταγμένα υπο-files και αποθηκεύουμε το αποτέλεσμα εναλλάξ στις input και output ταινίες μέχρι να διατάξουμε το αρχικό file μήκους N .

π.χ. input: A SORTING AND MERGING EXAMPLE

Αρχικά, χωρίζουμε το input-file σε 3-δες ($M=3$) στοιχείων και τα διατάσσουμε:

AOS	DMN	AEX	1 ^η
IRT	EGR	LMP	2 ^η
AGN	GIN	E	3 ^η

1^ο block

Επειτα, παίρνουμε τα πρώτα στοιχεία του 1^{ου} block της κάθε ταινίας και τα διατάσσουμε. Αυτό τοποθετώνται ανα 9-δες (διατεταγμένες) σε 3-output ταινίες:

AAG	INORST	1 ^η
DEGG	INNRR	2 ^η
AEE	LMPX	3 ^η

2^ο βήμα

Τελικά, παίρνουμε μια ταινία στην οποία διατάσσουμε ξανά τα στοιχεία των 3 ταινιών εξόδου (από το μικρότερο στο μεγαλύτερο).

AAADEEEEEGGG...

- Στο 1^ο βήμα έχουμε N/M blocks, όπου M είναι το μέγεθος της μνήμης και N είναι ο αριθμός των στοιχείων (γραμμάτων) σε κάθε block. (στο παράδειγμα έχουμε $M=3$ και $N=9$)
- Έστω k το μέγιστο πλήθος των βημάτων αναμίξης. Στο i -οστό βήμα παράγονται $p^i \cdot (\frac{N}{M})$ blocks. (όπου N είναι το συνολικό μήκος του file)
- Ο αλγόριθμος τελειώνει όταν έχουμε μόνο 1 block, δηλαδή $(\frac{N}{M}) \cdot \frac{1}{p^k} = 1 \Rightarrow k = \log_p(\frac{N}{M})$ (στο προηγούμενο παράδειγμα έχουμε $N=25$ και $M=3$, άρα $\lceil N/M \rceil = 9$ και τότε $k = \log_3 9 \Rightarrow k=2$ βήματα).

είναι το συνολικό $\lceil N/p \rceil (= \lceil 25/3 \rceil = 9)$

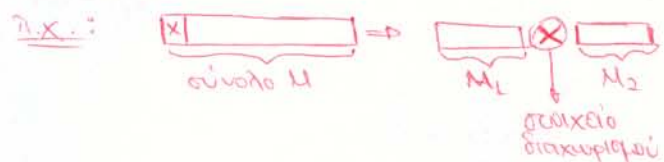
συνολικό N

ΓΡΑΜΜΙΚΟΣ MEDIAN ΑΝΤΙΟΡΙΘΜΟΣ

- Δίνεται μια ακολουθία από ανά δύο διαφορετικά στοιχεία και ένας αυθαίρετος i με $1 \leq i \leq n$ (η θέση του στοιχείου της ακολουθίας) και θέλουμε να βρούμε το i -οστό μικρότερο στοιχείο x_i τέτοιο ώστε να υπάρχουν $i-1$ στοιχεία x_k με $x_k < x_i$ και $n-i$ στοιχεία x_k με $x_k > x_i$.
- Αυτό το $n/2$ -οστό μικρότερο στοιχείο ονομάζεται median
- Μια πρώτη λύση είναι να διατάξουμε την ακολουθία και να απαριθμήσουμε έτσι ώστε να βρούμε το ζητούμενο στοιχείο. Κάτι τέτοιο όμως απαιτεί χρόνο $O(n \log n)$ στη χειρότερη και μέση περίπτωση.
- Ο αλγόριθμος find έχει γραμμικό χρόνο στη μέση περίπτωση ($O(n)$) και $O(n^2)$ στη χειρότερη

Η διαφορά του find με τον Quicksort στη μέση περίπτωση, αφού και οι δύο βασίζονται στην αρχή divide and conquer, είναι στο ότι ο Quicksort θέλει ταυτόχρονα και τα 2 υποπρόβλήματα, ενώ ο find

- ... ένα από αυτά. Σαν χειρότερη, όμως, περίπτωση και ο find δεν απορρίπτει τίποτα.
- Ο αλγόριθμος find διαχωρίζει το σύνολο M με βάση ένα οποιοδήποτε στοιχείο i του M σε δύο υποσύνολα, το M_1 και το M_2 . Αν $M_1 < i$, τότε το i -οστό μικρότερο στοιχείο δεν είναι στο υποσύνολο M_1 και είναι στο M_2 . Αν $M_1 = i$ τότε το στοιχείο x είναι το στοιχείο διαχωρισμού, δηλαδή το i -οστό μικρότερο στοιχείο. Αν $M_1 > i$ τότε το i -οστό μικρότερο στοιχείο θα είναι στο M_1 .



έχουμε την ακολουθία 5, 6, 4, 3, 6, 7
ως στοιχείο διαχωρισμού έχω το 5 = x
οπότε το M_1 θα αποτεθείται απ' τα στοιχεία που είναι μικρότερα του 5 και το M_2 απ' τα μεγαλύτερα.

$$\underbrace{(4, 3)}_{M_1} \quad (5) \quad \underbrace{(6, 6, 7)}_{M_2}$$

Ψάχνουμε το i -οστό μικρότερο στοιχείο της ακολουθίας. Αν π.χ. $i=4$, το 4^ο μικρότερο στοιχείο της ακολουθίας είναι το 6 (στην παρούσα περίπτωση το 4^ο μικρότερο στοιχείο ανήκει στο υποσύνολο M_2). Οπότε, τα M_1 και x πετούνται και ψάχνεται το M_2 . Φυσικά, προχωρώντας στον αλγόριθμο και καθλώντας τον αναδρομικά για το M_2 , θα έχουμε ως νέο x το 6 και τώρα θα ψάχνουμε το $i - |M_1| - 1$ στοιχείο, δηλαδή το $i - 2 - 1 = 4 - 3 = 1^{\circ}$ στοιχείο, το οποίο τυχαία να είναι ίδιο με το στοιχείο διαχωρισμού x . Οπότε το i -οστό (4^ο) στοιχείο που είναι μικρότερο στην ακολουθία είναι το 6.

- Η πολυπλοκότητα χειρότερης περίπτωσης για τον find είναι $O(n^2)$ όπως φαίνεται, αν πάρουμε $i=1$ και $|M_1| = |M| - 1$.

Προβλεπόμενες μέσες περιπτώσεις: υποθέτουμε ότι τα στοιχεία του M είναι από δύο διαφορετικά και κάθε μετέθεσή τους είναι ισοπρόβλητη. Θεωρώντας ότι το στοιχείο s που επιλέγεται ως στοιχείο διαχωρισμού επιλέγεται με πιθανότητα $p = \frac{1}{|M|}$ (όπου $|M|=n$)

Τότε: $T(n, i) \leq c \cdot n + \frac{1}{n} \left[\sum_{k=0}^{i-2} T(n-k-1, i-k-1) + \sum_{k=i}^{n-1} T(n-k, i) \right]$

($T(1)=c, c$ σταθ.)

$p = \frac{1}{|M|}$

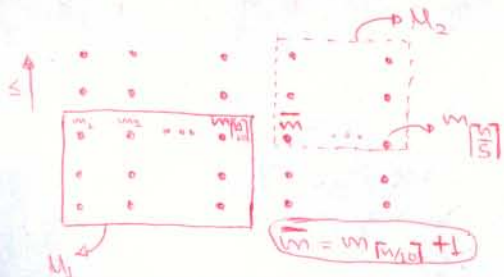
μέσος χρόνος για find($M_1, i - |M_1| - 1$) εφόσον $M_1 < i - 1$

μέσος χρόνος για find(M_2, i) εφόσον για $|M_2| > i - 1$ δηλαδή για find(M_2, i)

αρκεί $T(n) = \max_i T(n, i)$

Εμπειρικά αποδεικνύεται ότι $T(n) \leq 4cn$

- Για τον αλγόριθμο Select, αν έχουμε μεγάλο σύνολο M (με $n \geq 100$) τότε χωρίζουμε το input σε 5-δες. Έπειτα, διατάσσουμε τις πεντάδες και βρίσκουμε το median κάθε πεντάδας. Τέλος, κρατάμε όλους τους medians και βρίσκουμε το median των medians.



χωρίς βλάβη της γενικότητας θεωρούμε ότι $m_1, \dots, m_{\lfloor n/5 \rfloor} < m \leq m_{\lfloor n/5 \rfloor + 1}, \dots, m_{\lfloor n/5 \rfloor}$. Με βάση το m θα εφαρμόσαμε αναδρομικά τον select στο M_1 ή στο M_2 .

Είναι: $|M_1|, |M_2| \leq 3n/11$ και κάθε ερμηνεία έχει $\frac{3n}{10}$ στοιχεία.

Αν $T(n)$ είναι ορίσμος χρόνος για τον select τότε υπάρχουν σταθερές a, b

οι οποίες:

$$\begin{cases} T(n) \leq a \cdot n, & n \leq 100 \\ T(n) \leq T(\frac{3n}{10}) + T(\frac{7n}{10}) + b \cdot n, & n > 100 \end{cases}$$

Εμπειρικά αποδεικνύεται και στις 2 περιπτώσεις ότι $T(n) \leq cn$
Άρα ο select εκτελείται σε γραμμικό χρόνο.

διαχωρίζουμε το M σε $\lfloor n/5 \rfloor$ υποσύνολα με 5 στοιχεία το καθένα.

5.4. ΠΟΤΗΤΕΣ ΔΟΜΕΣ ΔΕΔΟΜΕΝΩΝ

Οι συνιστώσες δομής δεδομένων αποθηκεύουν ένα σύνολο M , όπου $M \in U$ και U σύνολο με $M = \{x_1, \dots, x_n\}$, και χρειάζονται χώρο $g \cdot n + c_2$ όπου g, c_2 σταθερές. Το $c_1 = 1$ και το c_2 είναι μικρό.

• Αν το σύνολο M είναι διατεταγμένο και αποθηκευμένο σε έναν πίνακα $S = S[1 \dots n]$ με $S[i] = x_i$ και $x_1 < x_2 < \dots < x_i < \dots < x_n$, τότε ο παρακάτω αλγόριθμος υλοποιεί την πράξη search(x) κατά την οποία η πληροφορία των στοιχείων αντλείται μόνο με συγκρίσεις:

```

left = 1, right = n
next = τυχαίος αριθμός από [left...right];
while x ≠ S[next] and right > left do
    if x < S[next] then right = next - 1
    else left = next + 1
next = τυχαίος αριθμός από [left...right];
if x = S[next] then επανείσοδος
else απουσία
    
```

Αν ισχύει: • next = left \rightarrow έχουμε linear search
 • next = $\lfloor \frac{right+left}{2} \rfloor \rightarrow$ έχουμε binary search
 • next = $\lfloor \frac{x - S[left]}{S[right] - S[left]} \cdot (right - left) \rfloor + left \rightarrow$ έχουμε interpolation search
 ↓
 ισχύει για uniform κατανομή, δηλαδή όλα τα στοιχεία είναι ισοπλάσια για να εμφανισθούν

1

• Binary interpolation search: Για να βρούμε ένα στοιχείο a στο διατεταγμένο σύνολο $M = \{x_1, \dots, x_n\}$ με $x_1 < x_2 < \dots < x_n$ κάνουμε τα εξής:

- Αρχικά, εφαρμόζουμε interpolation search για το a , καθορίζοντας εύρος του που μπορεί να βρεθεί.
- Ψάχνουμε σε $\lfloor \sqrt{n} \rfloor$ αποστάσεις ~~από το εύρος~~ κινώντας αλγόριθμο μεγέθους \sqrt{n} έτσι ώστε να βρούμε σε ένα στοιχείο $z = x_k > a$.
- Εφαρμόζουμε την ίδια διαδικασία στο υποφίλε μεγέθους $\lfloor \sqrt{n} \rfloor$ στο οποίο περιέχεται το a .

• Πομπανωδότητα χειρότερης περίπτωσης: το μέγιστο πλήθος των συγκρίσεων είναι $\sum_{i=1}^{\infty} \frac{1}{2^i} = 1$
 $= \sqrt{n} + \sqrt{n} + \sqrt{n} + \dots = O(\sqrt{n})$

• Πομπανωδότητα μέσης περίπτωσης: θεωρούμε ότι τα στοιχεία x_i που επιλέγονται από το σύνολο υπακούουν σε uniform κατανομή. Έστω τότε G ο μέσος αριθμός των συγκρίσεων που απαιτείται για να βρούμε το κατάλληλο υποφίλε μεγέθους \sqrt{n} και P_i η πιθανότητα του να χρειάζονται τουλάχιστον i συγκρίσεις για να βρεθεί αυτό το υποφίλε. Τότε:

$$G = \sum_{i \geq 1} i(P_i - P_{i+1}) = \sum_{i \geq 1} P_i \text{ με } P_1 = P_2 = 1, \text{ αφού πάντα εκτελούνται τουλάχιστον 2 συγκρίσεις}$$

Ισχύει: $P_i \leq \text{Prob}\{| \theta_{\text{επιλεγ}} - \alpha - n\epsilon | > (i-2)\sqrt{n}\}$ για $i \geq 3$

Chebyshev inequality: $\text{Prob}(|x - \mu| \geq \epsilon) \leq \sigma^2 / \epsilon^2$, όπου ϵ σταθερά, μη μέση τιμή και σ η απόκλιση.
 Αν x είναι τυχαία μεταβλητή που παίρνει τιμές μικρότερες του α , τότε θα δείξω ότι $x = \theta_{\text{επιλεγ}}$ του α , $\mu = \text{next}$ και $\sigma^2 / \epsilon^2 = (i-2)\sqrt{n}$

Κάνω πειράματα μεγέθους \sqrt{n} στο $[x_0, x_{n+1}]$. Αν κάνω i τέτοια πειράματα τότε η πιθανότητα P_i να είναι $\theta_j = \theta(x_j)$. Η πιθανότητα να επιλέξω ένα στοιχείο μικρότερο του α είναι:

$$P = \frac{\alpha - x_0}{x_{n+1} - x_0} \cdot \text{Τότε } \theta_j = \binom{n}{i} \cdot p^i \cdot (1-p)^{n-i}, \text{ δηλαδή έχω binomial κατανομή με } \mu = p \cdot n \text{ και } \sigma^2 = p(1-p)n \leq \frac{n}{4}$$

Αρα $\mu = \text{next}$ και θα είναι:

$$P_i \leq \text{Prob}\{|x - \mu| \geq (i-2)\sqrt{n}\} \leq \frac{p(1-p)}{((i-2)\sqrt{n})^2} = \frac{p(1-p)}{(i-2)^2} \leq \frac{1}{4(i-2)^2}$$

$$\sigma^2 = p(1-p)n \leq \frac{n}{4} \text{ διασπορά (για } p = \frac{1}{2} \Rightarrow \sigma^2 = \frac{n}{4})$$

$$\text{Οπότε: } G \leq \frac{1+1}{2} + \sum_{i \geq 3} \frac{1}{4(i-2)^2} \approx 2.4$$

Αν, λοιπόν, είναι $\bar{T}(n)$ μέσο πλήθος των συγκρίσεων τότε $\bar{T}(n) \leq G + \bar{T}(\sqrt{n})$, για $n \geq 3$ $\frac{\bar{T}(n) \leq 1}{\bar{T}(n) \leq 2} \Rightarrow \bar{T}(n) \leq 2 + 2.4 \log \log n$

$\Rightarrow \bar{T}(n) = O(\log \log n)$ Άρα ο μέσος χρόνος για το BTS είναι $O(\log \log n)$

• Για να διορθώσουμε την πολυπλοκότητα χειρότερης περίπτωσης του BTS εφαρμόσαμε το ευθείως και διίως ψάξιμο. Αντί δηλαδή να εφαρμόσουμε γραμμικό ψάξιμο για να βρούμε ένα τέτοιο ώστε $i \cdot n$ βήματα να έχουμε βρεί το υποφίλε που βρίσκεται το a , καίνουμε ευθείως με βήματα $2^0, 2^1, \dots, 2^j$ ώστε να βρούμε ένα j τέτοιο ώστε με $2^j \cdot n$ βήματα να βρούμε το υποφίλε το a . (ιχίει $j < \log n$)

• Έστω ότι βρίσκουμε ένα υποφίλε για το οποίο ιχίει $S[\text{next} + 2^{j-1} \cdot n] < a \leq S[\text{next} + 2^j \cdot n]$. Τότε

σ' αυτό το υποφίλε εφαρμόσουμε διίως ψάξιμο. Δηλαδή:



Το υποφίλε έχει μέγεθος $2^j \cdot n - 2^{j-1} \cdot n = 2^{j-1} \cdot n$. Μπορώ, λοιπόν, να το χωρίσω σε n υπο-ομάδες που να έχω 2^{j-1} υπο-ομάδες. Απώ, δηλαδή, βρώ το διάστημα $2^{j-1} \cdot n$ εφαρμόζω μέσα σε αυτό το γραμμικό ψάξιμο για να βρώ το n διάστημα στο οποίο περιέχεται το στοιχείο που. Έπειτα, σ' αυτό το διάστημα μεγέθους n θα εφαρμόσω πάλι την ίδια μέθοδο, δηλαδή τα ευθείως βήματα και το γραμμικό ψάξιμο.

Το διίως ψάξιμο έμειται στο γεγονός ότι πόλλες φορές το υποφίλε τότε σ' αυτό ψάχνουμε διίως, δηλαδή ψάχνουμε διίως μετάξυ των κομμάτιών $2^{j-1} \cdot n, (2^{j-1} + 1) \cdot n, \dots, 2^j \cdot n$. Με άλλα λόγια τα διίως σφάλματα που δεν χρειάζονται τα τεταίω. (αν π.χ. έχω 8 υπο-ομάδες ($j=3$) τις χωρίζω σε 4 και 4 και κατόω που είναι το a ...)

• Το ευθείως ψάξιμο χρειάζεται χρόνο $O(\log n)$ και το διίως ψάξιμο χρειάζεται $\log 2^{j-1} = O(j-1) = O(\log n)$

Άρα ο χειρότερος χρόνος είναι $O(\log n)$ διότι $\frac{1}{2} \log n$ ο συνολικός χειρότερος χρόνος είναι:

$$\log n^{1/2} + \log n^{1/4} + \dots = (1/2 + 1/4 + \dots) \cdot \log n = O(\log n)$$



ΕΚΤΕΝΕΙΣ ΔΟΜΕΣ ΔΕΔΟΜΕΝΩΝ

Οι εκτενείς δομές δεδομένων αποθηκεύουν ένα σύνολο M , με $|M|=n$ και χρειάζεται χώρο $c_1 \cdot n + c_2$ όπου c_1, c_2 σταθερές με $c_1 \neq 1$ και c_1, c_2 μικρές, κυρίως $c_1 \in \{3, 4, 5\}$

Διακρίνονται στις απλές που υλοποιούνται με χρήση δεικτών και στις υβριδικές που χρησιμοποιούν και αττάμς.

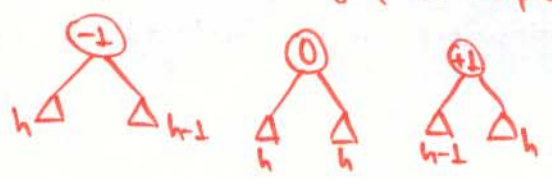
• Ισοζυγισμένα δένδρα: Ένα δένδρο εύρεσης T για το σύνολο S με $|S|=n$, ονομάζεται ισοζυγισμένο όταν ισχύει: $h(T) = c_1 \log n + c_2$, c_1, c_2 σταθερές και $h(T)$ το ύψος του T .

⊗ Τα ισοζυγισμένα δένδρα ορίζονται βάσει ενός κριτηρίου ζυγισής και εκτελούν το ψάξιμο, ένθεση, απόσβεση σε $O(\log n)$ χρόνο. (συν. χ.π.).

α) Υποζυγισμένα δένδρα: ισοζυγισουμε βάσει του ύψους των υπο δένδρων κάθε κόμβου.

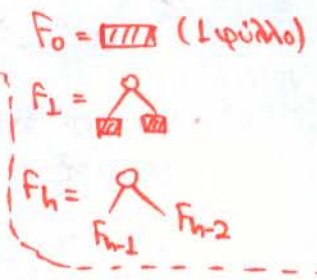
1) AVL-δένδρο

- Έστω $h_b(u)$ η υποζύγισση του κόμβου u , με $h_b(u) = \text{ύψος}(\text{τσον}(u)) - \text{ύψος}(\text{λεσον}(u))$ όπου τσον ο δεξιός γιος και λεσον ο αριστερός
- Ένα δυϊκό δένδρο T λέγεται AVL όταν για κάθε κόμβο $u \in T$ ισχύει $|h_b(u)| \leq 1$
- Υπάρχουν 3 κατηγορίες κόμβων ανάλογα με την υποζύγισση:



} Ο κόμβος με υποζύγισση 0 είναι ζυγισμένος και οι άλλοι μη ζυγισμένοι.

- Τα δένδρα - Fibonacci ορίζονται ως εξής (F_0, F_1, F_2, \dots):
Το F_h δένδρο έχει ύψος h



⊗ Το F_h περιέχει ακριβώς $\frac{a^{h+2} - b^{h+2}}{\sqrt{5}}$ φύλλα, με $a = \frac{1+\sqrt{5}}{2}$, $b = \frac{-1+\sqrt{5}}{2}$ και ο αριθμός φύλλων δίνεται απ' την ακολουθία $|F_h| = |F_{h-1}| + |F_{h-2}|$ με $F_0 = 1, F_1 = 2$.

⊗ Για AVL με n φύλλα και ύψος h ισχύει: $|F_h| \leq n \leq$ πλήρες δυϊκό με ύψος $h \Rightarrow$ (πλέον ισορροπημένο AVL)
 $\Rightarrow \frac{a^{h+2} - b^{h+2}}{\sqrt{5}} \leq n \leq 2^h$. Είναι: $n \leq 2^h \Rightarrow \log n \leq h$

Όμως: $\frac{a^{h+2} - b^{h+2}}{\sqrt{5}} \leq n \xrightarrow{b \ll 1} n \geq \frac{a^{h+2} - 1}{\sqrt{5}} \Rightarrow a^{h+2} \leq \sqrt{5}n + 1 \Rightarrow a^{h+2} \leq \sqrt{5}(n+1) \Rightarrow$
 $\Rightarrow \dots \Rightarrow h \leq 1.44 \log(n+1) - 0.3277$
(υφίσταται από λογάριθμο).

\Rightarrow Το ψάξιμο στο AVL κοστίζει $O(\log n)$

⊗ Οι εξαιρετικότητες καταστάσεις για το AVL είναι το F_h και το πλήρες δυϊκό δένδρο (όλα με ύψος h)

• Ενθεση-απόσβεση στο AVL:

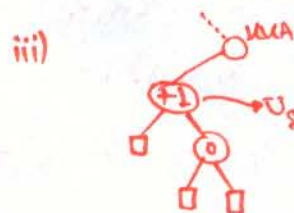
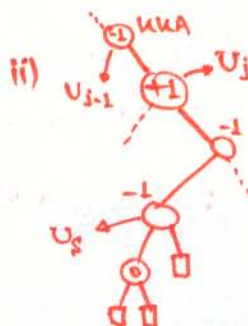
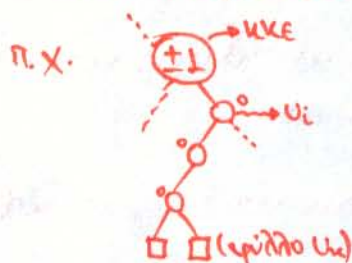
→ Έστω μονοπάτι $p = (u_0, u_1, \dots, u_k)$ απ' τη ρίζα u_0 έως το φύλλο u_k . Αν για τον κόμβο $u_i \in p$ ισχύει $hb(u_i) = hb(u_{i+1}) = \dots = hb(u_{k-1}) = 0$ τότε ο κόμβος u_{i-1} ονομάζεται κρίσιμος κόμβος ενθέσεως (ΚΚΕ).

→ Έστω u_s ένας κόμβος στο μονοπάτι ελεύθερος με $h(u_s) = 2$. Τότε ο ΚΚΑ είναι:

i) ο u_s όταν $hb(u_s) = 0$

ii) ο u_{j-1} όπου $j \leq t < s$ για $hb(u_s) = \begin{cases} +1, & \text{όταν } \tau_{\text{son}(u_s)} \text{ στο μονοπάτι ελεύθερος} \\ -1, & \dots \text{son}(u_s) \end{cases}$ " " " "

iii) ο $\pi(u_s)$ όταν δεν υπάρχει t που να ικανοποιεί ανυ προηγούμενη σχέση.



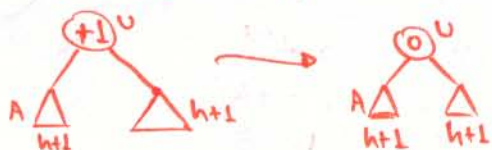
⊗ Μετά από κάθε ενθεση ή απόσβεση το ύψος του υποδένδρου με ρίζα ΚΚΕ ή ΚΚΑ μεταβάλλεται κατά $+1$ ή -1 (απόσβεση). Πάνω στο κρίσιμο μονοπάτι (u_i, \dots, u_k) εκτελούνται οι αλλαγές δόχους. Γενικά, οι ενθέσεις ή αποσβέσεις αλλάζουν τα ύψη του υποδένδρου και έτσι μπορεί να παρουσιασθεί διαφορά $+2$ ή -2 και γ'αυτ' εκτελούνται επαναληψιμικές πράξεις.

⇒ Ενθεση

Έστω u ο ΚΚΕ και χωρίς βλάβη της γενικότητας είναι $hb(u) = +1$ και έχουμε το δένδρο:



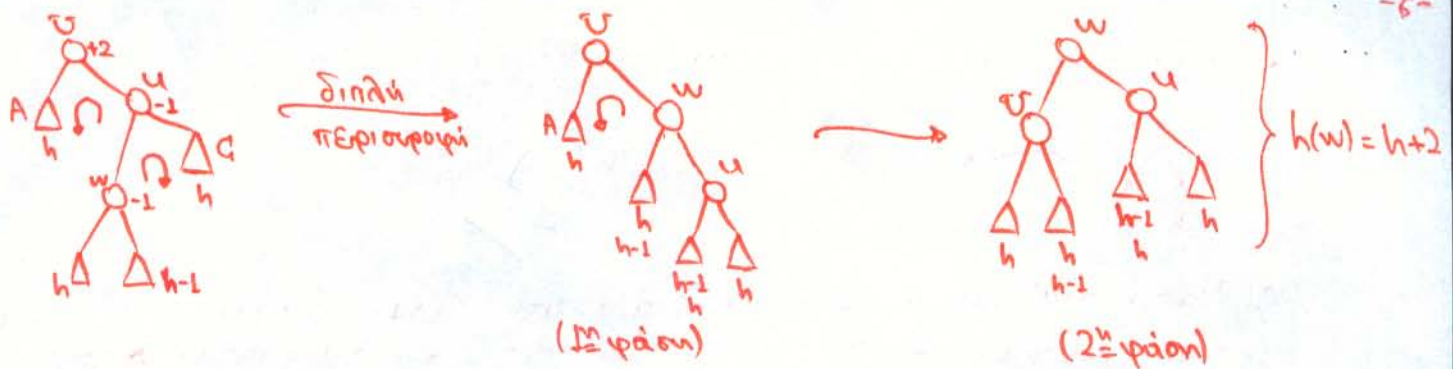
• Κατόπιν ενθέσεως στο A έχουμε απορρόφηση (absorption): το ύψος του A γίνεται $h+1$, οπότε θέτουμε $hb(u) = 0$



• Κατόπιν ενθέσεως στο G έχουμε απλή περιστροφή (rotation): το ύψος του G γίνεται $h+1$ οπότε ο u έχει βάρος $+2$ ($+1+1$) και ο u $+1$



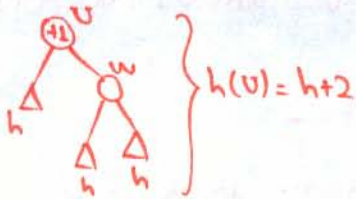
• Κατόπιν ενθέσεως στο B έχουμε διπλή περιστροφή (double rotation): θεωρούμε ότι η ενθεση γίνεται στο αριστερό υποδένδρο του B, δηλ. $h(l[B]) = (h-1)+1 = h$ και το δεξί υποδένδρο του B έχει ύψος $h-1$. Οπότε:



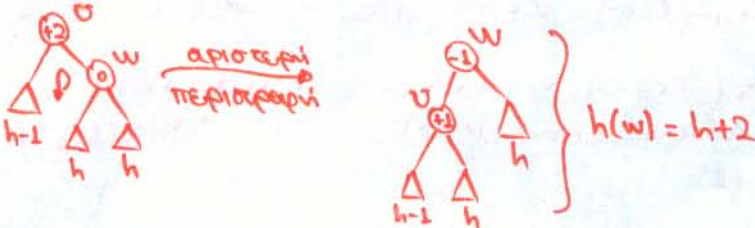
* Μετά από 1 ενθεση χρειάζονται $O(1)$ rotations ή double rotations
(το αλγόριθμος των υψώνων μονοπάτιών μιας ακολουθίας η ενθέσεων είναι 2,618n → κατανεμημένα 2,618).

⇒ Απόσβεση

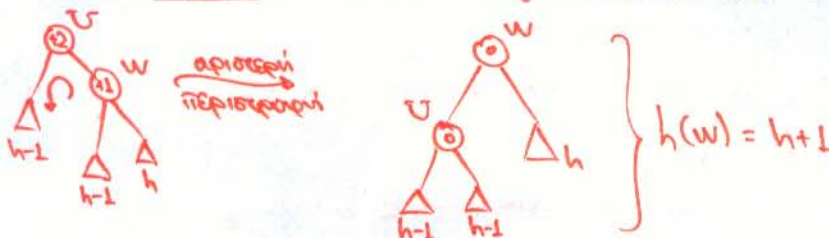
Έστω u ο ΚΚΑ και χωρίς βλάβη της γενικότητας είναι $hb(u) = +1$. Θα μελετήσουμε την απόσβεση στο αριστερό υποδένδρο με αλλαγή ύψους αυτού του υποδένδρου.



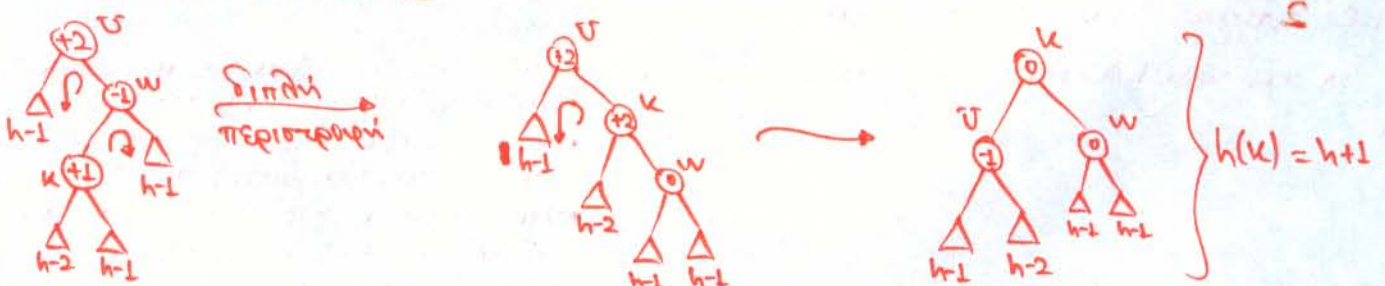
- Τελική περιστροφή (terminated rotation): αφαιρώντας απ' το αριστερό υποδένδρο, το ύψος του γίνεται $h-1$ οπότε $hb(u) = +2$ και $hb(w) = 0$:



- Συνεχισμένη περιστροφή (propagated rotation): αφαιρώντας απ' το αριστερό υποδένδρο του w , το ύψος του γίνεται $h-1$ και τότε $hb(u) = +2$ και $hb(w) = +1$:

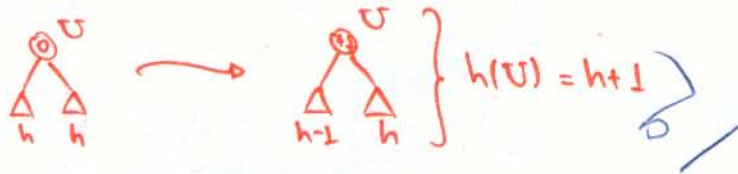


- Συνεχισμένη διπλοπεριστροφή (propagated double rotation): αφαιρώντας απ' το αριστερό υποδένδρο του u , το ύψος του γίνεται $h-2$ και τότε $hb(u) = +1$, $hb(w) = -1$ και $hb(v) = (h+1) - (h-1) = +2$



* Για $hb(u) = 0$ τότε στο τέλος έχω $hb(u) = 0$, $hb(v) = 0$, $hb(w) = 0$
Για $hb(u) = -1$ " " " " $hb(u) = 0$, $hb(v) = 0$, $hb(w) = +1$

• Απορρόφηση (absorption): όταν το $h_b(u) = 0$, δηλαδή:



⊛ Μετά από μια απόσβεση χρειάζονται $O(\log n)$ rotations ή double rotations. (κατανεμημένο κόστος για 1 απόσβεση σε μια ακολουθία n αποσβέσεων είναι 1.618)

2) Red-black δένδρο ή BB-δένδρο

- κάθε κόμβος του δένδρου έχει ακριβώς 2 παιδιά (πλήρες δυϊκό δένδρο \Rightarrow αδύνατες οι αλυσίδες κόμβων)
- τα στοιχεία αποθηκεύονται στα φύλλα
- σε κάθε κόμβο u αντιστοιχεί ένας ανεξάρτητος βαθμός (rank) που συμβολίζεται με $\tau(u)$
- αν το u είναι φύλλο τότε $\tau(u) = 0$ και $\tau(p(u)) = 1$
- ισχύει $\tau(u) \leq \tau(p(u)) \leq \tau(u) + 1$ (το $\tau(u)$ αυξάνεται κατά 1 ή μένει αμετάβλητο από πάνω προς τα κάτω)
- ισχύει $\tau(u) < \tau(p(p(u)))$

⊛ Αν έχουμε $\text{rank}(p(u)) - \text{rank}(u) = \text{bit}(u)$ τότε το $\text{bit}(u) \in \{0, 1\}$

Στο red-black δένδρο έχουμε τις απλές πράξεις επανατάξης στις οποίες έχουμε αλλαγή χρώματος (promotions, demotions) και τις σύνθετες πράξεις ή δομικές (rotations ή d-rotations)

- Κάθε ένθεση ή απόσβεση απαιτεί $O(\log n)$ απλές πράξεις (promotions-demotions) και $O(1)$ δομικές (rotations)
- Μια ακολουθία από $O(n)$ ενθέσεις ή αποσβέσεις σε ένα αρχικά άδειο BB δένδρο απαιτεί $O(n)$ promotions ή demotions.

- ⊙ Ένας κόμβος σε ύψος t θα λάβει μέρος σε μια επανατάξη με πιθανότητα $(\frac{1}{2})^t$
 - ⊙ Είναι δυνατή η top-down επανατάξη
 - Ένα bit είναι αναγκαίο για την αποθήκευση της τύξης.
- \rightarrow είναι πολύ βασικές ιδιότητες σε εφαρμογές παραλληλισμού και πολυδιάστατων δομών δεδομένων

Όσον αφορά το χρώμα των αιχμών, αυτές θα είναι κόκκινες (τιμή 0) αν οι διαφορές των ranks των 2 κόμβων της αιχμής είναι μηδενικές και μαύρες (τιμή 1) αν οι διαφορές των ranks είναι $\neq 0 (=1)$. (αφού $\tau(u) < \tau(p(p(u)))$ δεν μπορούμε να έχουμε 2 αιχμές κόκκινες στη σειρά)

Για το χρωματισμό των κόμβων έχουμε:

- κάθε κόμβος είναι κόκκινος ή μαύρος
- κάθε φύλλο είναι μαύρο
- αν ένας κόμβος είναι κόκκινος τότε και τα 2 παιδιά του θα είναι μαύρα.
- κάθε αιχμή μονοπάτι από ένα κόμβο προς ένα φύλλο που είναι απόγονος του περιέχει τον ίδιο αριθμό από μαύρους κόμβους

β) Βαροζυγισμένα δένδρα: 100 ζυγίζουμε βάσει του πλήθους των κόμβων των υποδένδρων ενός κόμβου.

1) BB[a]-δένδρο

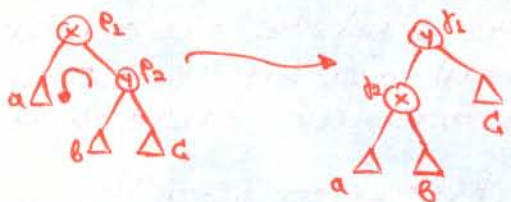
- Έστω T ένα δυϊκό δένδρο με αριστερό υποδένδρο T_L και δεξί υποδένδρο T_R . Τότε, ορίζουμε ως $p(T) = \frac{|T_L|}{|T|}$ τη βαροζύγηση της ρίζας του T , όπου $|T|$ είναι ο αριθμός των φύλλων του δένδρου.
- Ένα δένδρο T έχει βαροζύγηση a όταν \forall υποδένδρο T' του T ισχύει: $a \leq p(T') \leq 1-a$
- Ένα $BB[a]$ -δένδρο είναι το σύνολο όλων των δένδρων T περιφορμένων βαροζύγησης a
- Όταν ισχύει $a \leq p(T')$ έχουμε ποσά στοιχεία στα αριστερά (άρα βαδίζουμε δεξιά)
- " " " $p(T') \leq 1-a$ " " " δεξιά (" " αριστερά)

(*) Έστω δένδρο T το οποίο είναι $BB[a]$ -δένδρο. Τότε: $h(T) \leq 1 + \frac{\log(n+1)-1}{\log \frac{1}{1-a}} = \alpha \log n$
(το ύψος του T φράσσεται από λογάριθμο)

Όταν οι ενθέσεις ή οι αποσβέσεις επηρεάζουν τις βαροζυγίες των κόμβων και τις βγαίνουν έξω απ' το διάστημα $[a, 1-a]$ τότε μπορούμε να επαναφέρουμε τη ζύγηση στο επιτρεπτό διάστημα με απλές ή διπλές περιστροφές.

→ απλή περιστροφή:

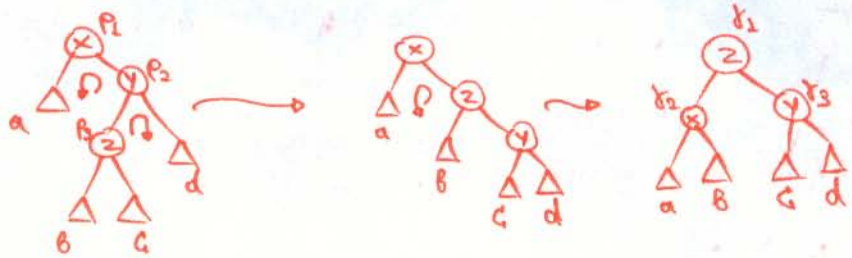
$p_1, p_2, \delta_1, \delta_2$ είναι οι βαροζυγίες των αντίστοιχων κόμβων



όπου $\delta_1 = p_1 + (1-p_1)p_2$
 $\delta_2 = \frac{p_2}{p_1 + (1-p_1)p_2} = \frac{p_2}{\delta_1}$

→ διπλή περιστροφή:

$p_1, p_2, p_3, \delta_1, \delta_2, \delta_3$ είναι οι βαροζυγίες των αντίστοιχων κόμβων



όπου $\delta_1 = p_1 + (1-p_1)p_2p_3$, $\delta_2 = \frac{p_2}{p_1 + (1-p_1)p_3} = \frac{p_2}{\delta_1}$, $\delta_3 = \frac{p_3(1-p_3)}{1-p_2p_3}$

(*) Για όλα τα $a \in (\frac{1}{4}, 1 - \frac{\sqrt{5}}{2}]$ υπάρχουν σταθερές $d \in [a, 1-a]$ και $\delta \geq 0$ τέτοιες ώστε για ένα δυϊκό δένδρο T με υποδένδρα T_L και T_R (τα οποία είναι $BB[a]$ -δένδρα) να έχουμε $p(T) = \frac{|T_L|}{n}$ και:

$p(T) = \frac{|T_L|}{|T|} < a$ και

$\frac{|T_L|}{|T|-1} \geq a$ (δηλαδή έγινε ένθεση στο δεξί υποδένδρο)

$\frac{|T_L|+1}{|T|+1} \geq a$ (" " αποσβέση στο αριστερό ")

- τότε θα ισχύει
- i) αν $p(T_L) \leq d$ τότε μια απλή περιστροφή επαναφέρει τη ζύγηση και $\delta_1, \delta_2 \in [(1+\delta)a, 1-(1+\delta)a]$
 - ii) αν $p(T_L) > d$ τότε μια διπλή περιστροφή επαναφέρει τη ζύγηση και $\delta_1, \delta_2, \delta_3 \in [(1+\delta)a, 1-(1+\delta)a]$

Γενικά, μετά την ένθεση ή την αποσβέση εφαρμόζουμε επαναφορτισμένες πράξεις στο μονοπάτι απ' το φύλλο προς τη ρίζα οι οποίες στη χειρότερη περίπτωση είναι $O(\log n)$. Άρα, στη χ.π. έχουμε για φάξιμο, ένθεση και αποσβέση $O(\log n)$ πολυπλοκότητα.
Η ιδιότητα βάρεως εγγυάται μια ικανοποιητική μακροεξέταση των αριθμών επαναζυγιστικών πράξεων και έτσι μπορεί να γίνει καταμερίσιμος (αμορτίζαται)

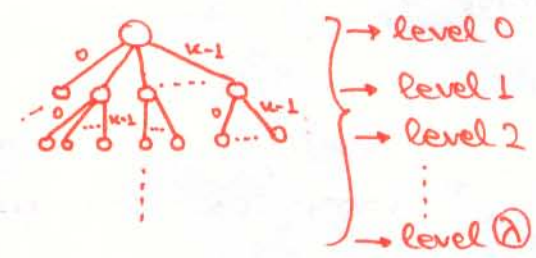
Υβριδικές δομές δεδομένων

1) Τρίες (ψηφιακά δένδρα): προτάθηκε ως ένας απλός τρόπος δομής ενός file χρησιμοποιώντας την ψηφιακή αναπαράσταση των στοιχείων του.

Ορίζουμε ένα σύμπαν $U = \{0, \dots, u-1\}^\lambda$ το οποίο αποτελείται από strings μήκους λ ψηφίων τα οποία δομούνται πάνω σε ένα αλφάβητο με κ στοιχεία. ($κ \geq 2$)

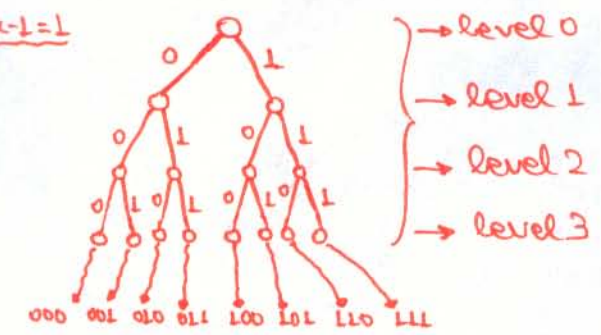
- Ένα σύνολο $S \subseteq U$ αναπαρίσεται ως ένα κ-ικό δένδρο που περιέχει όλα τα πρόθεμα (prefixes) των S (των στοιχείων του).
- Μια προφανής υλοποίηση ενός τριε είναι να χρησιμοποιήσουμε για κάθε εσωτερικό κόμβο του δένδρου ένα ατταχ μήκους κ.

Αναπαράσταση του συνόλου $S \subseteq U$:



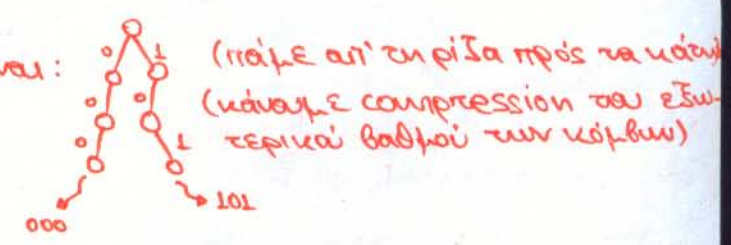
Στην περίπτωση της αναπαράστασης του τριε που αποθηκεύει όλα τα στοιχεία του σόμπαντος U θα είναι: $\lambda = \log_{\kappa} |U|$ όπου $|U|$ είναι το μέγεθος του σόμπαντος.

• Θεωρούμε το σύμπαν $U = \{0, 1\}^3$, όπου δηλαδή $κ=2$ στοιχεία (0,1) και $\lambda=3$, δηλαδή έχουμε strings των 3 ψηφίων. Τότε το μέγεθος του σόμπαντος U είναι $|U|=8$, αφού σ'αυτό αποθηκεύονται (περιέχονται) τα strings 000, 001, 010, 011, 100, 101, 110, 111. Θα αναπαράσσουμε το πλήρες τριε που αποθηκεύει όλα τα στοιχεία του U . Είναι:



Βλέπουμε ότι $h(\text{τριε}) = 3$ και το τελευταίο επίπεδο (level λ) είναι το 3. Οπότε: $\lambda = \log_2 |U| = \log_2 8 = 3$ (σωστό).
Ο αριθμός φύλλων του τριε είναι όσο και το μέγεθος του U . ($κ^\lambda = 2^3 = 8$)

Για το σύνολο $S = \{000, 101\} \subseteq U$ τα νέο τριε είναι:



Παρατηρούμε ότι θέλουμε χρόνο $O(1)$ για ψάξιμο.

Γενικά, ο χώρος που χρησιμοποιεί το τριε είναι υπερβολικός: $O(n \cdot \lambda \cdot \kappa)$, όπου λ είναι το ύψος του τριε, κ είναι τα στοιχεία του αλφάβητου του σόμπαντος και n είναι το μέγεθος του συνόλου S ($|S|=n$). Οπότε, έχουμε:

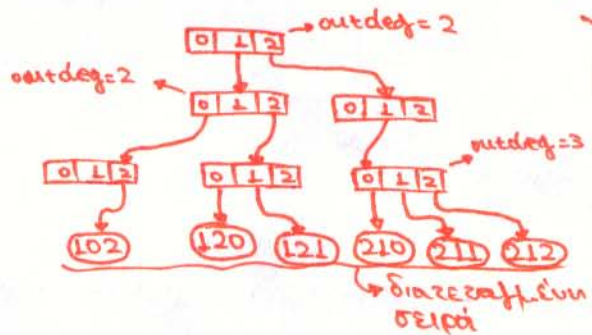
$$\text{χώρος} = n \cdot \lambda \cdot \kappa = n \cdot \log_{\kappa} |U| \cdot \kappa \Rightarrow \boxed{\text{χώρος} = n \kappa \cdot \frac{\log |U|}{\log \kappa}}$$

$$\text{χρόνος} = O(\log_{\kappa} |U|) = O\left(\frac{\log |U|}{\log \kappa}\right) \Rightarrow \boxed{\text{χρόνος} = O\left(\frac{\log |U|}{\log \kappa}\right)}$$

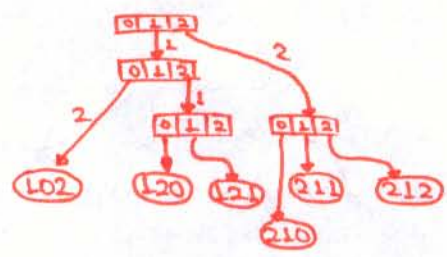
και ζώνεται όταν το ~~ατταχ~~ τότε έχουμε ελάχιστο χρόνο και μέγιστο χώρο. (όταν $\kappa \uparrow$ αντίθετα)

* ο χρόνος δεν εξαρτάται απ' το $n = |S|$ αλλά απ' το $|U|$, στη μέση περίπτωση όμως αυτός είναι $O(\log_{\kappa} n)$

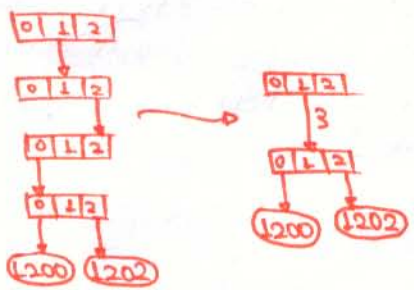
8) Για το σύνολο $S = \{121, 102, 211, 120, 210, 212\}$ χρησιμοποιήστε ένα αττάχ 3 θέσεων σε κάθε εσωτερικό κόμβο και παίρναμε το δένδρο:



Για να ελαττώσουμε το χώρο σε $O(n \cdot k)$ αποθηκεύουμε μόνο τους εσωτερικούς κόμβους που έχουν $outdeg \geq 2$. Τότε θα έχουμε το επόμενο σχήμα:



π.χ.



15

2) Interpretation search tree

- Ένα IST είναι ένα πολυδιαυλαδισμένο δένδρο όπου ο βαθμός διαυλάδωσης κάθε κόμβου εξαρτάται απ' το πλήθος των στοιχείων που αποθηκεύει στο υπόδενδρο του. Αντίθετα, ο βαθμός της ρίζας είναι $\Theta(n)$. Η ρίζα διαίρει το file σε $\Theta(\sqrt{n})$ υποfiles διαφορετικοί μεταξύ τους μήκους.
- Ένα Ισοκύβιο IST είναι ένα IST στο οποίο όλα τα υπο-files έχουν περίπου ίσο μέγεθος και επομένως οι γιοι της ρίζας έχουν βαθμό $\Theta(\sqrt{n})$. Τα Ισοκύβια IST έχουν βάθος $O(\log \log n)$.
- Ο αλγόριθμος ψαξίματος για το IST χρησιμοποιεί interpretation search σε κάθε κόμβο του IST για να βρει το υπο-file όπου θα συνεχιστεί το ψάξιμο. Γι' αυτό οι κόμβοι υλοποιούνται ως αττάχ. Σε κάθε αττάχ προσερχόμαστε και μια προσέγγιση της αντίστροφης συνάρτησης κατανομής, η οποία μας επιτρέπει να κάνουμε τέλει interpolate σε ένα δείγμα μεγέθους n^a , όπου $1/2 \leq a < 1$ είναι μια παράμετρος.
- Για να αποθηκευθούν n στοιχεία ο χώρος που χρησιμοποιείται είναι γραμμικός: $O(n)$.
- Για ψάξιμο σε ένα IST μεγέθους n , προφανώς χρειαζόμαστε $O(\sqrt{n})$ μονάδες χρόνου για το ψάξιμο στο αττάχ της ρίζας. Το επόμενο υποfile που θα ψαχθεί θα έχει μέγεθος το πολύ $n/2$ και συνεπώς θα ισχύει: $T(n) \leq O(\sqrt{n}) + T(n/2) \Rightarrow T(n) = O(\sqrt{n})$, όπου $T(n)$ ο χειρότερος χρόνος. Αυτός ο χρόνος βελτιώνεται σε $T(n) = O(\log^2 n)$ εφαρμόζοντας ευθείο και δυϊκό ψάξιμο.
- Για n ενθέσεις ή αποβρώσεις ο χρόνος χειρότερης περίπτωσης είναι $O(n \log n)$ και το αντίστοιχο αναμενόμενο κόστος (μέσος όρος) είναι $O(\frac{n \log n}{n}) = O(\log n)$. Ο χρόνος μέσης περίπτωσης με μορφή αναμενόμενου κόστους είναι $O(\log \log n)$.
- Σε κάθε κόμβο u του IST προσερχόμαστε έναν αριθμητικό $g(u)$. Αυτός αναριθμεί το πλήθος των ενθέσεων ή διαγραφών που γίνονται στο υπόδενδρο T_u του u . Αν αναμετασκευάσουμε το T_u του κόμβου u τότε δίνουμε αρχική τιμή 0 σε όλους τους μετρητές των κόμβων του T_u . Ένας αριθμητικός $g(u)$ υπερχηλιδεί (overflows) όταν η τιμή του ξεπεράσει το $n/4$, όπου n είναι το μέγεθος (size) του κόμβου u τη στιγμή που αναμετασκευάζουμε το δένδρο (δηλ. όταν $g(u) = 0$ την τελευταία φορά). Όταν γίνει overflow τότε γίνεται και η αντίστροφη αναμετασκευή.
- Ορίζουμε ως $w(u)$ το βάρος ενός κόμβου που είναι το πλήθος των στοιχείων που είναι αποθηκευμένα στο T_u . Το μέγεθος του u είναι το πλήθος των στοιχείων του T_u τα οποία είναι μη-μαρκιαρισμένα (μαρκιαρισμένο λέγεται ένα στοιχείο όταν αυτό πρόκειται να αποβρωθεί, δεν απομακρύνεται όμως παρά μόνο τότε που θα γίνει η αναμετασκευή του υποδένδρου που το περιέχει).

- Μια ένθεση στο IST γίνεται ως εξής: βρίσκουμε τη θέση που πρέπει να βρισκόταν το x και το εμβόλιμο με x αντιστοίχως. Έπειτα, αντιστρέφουμε τους μετρητές όλων των κόμβων v που βρίσκονται στο μονοπάτι απ' το x στη ρίζα. Βρίσκουμε τον υψηλότερο τέτοιο κόμβο, ώστε ο μετρητής του να υπερβεί, κι αν χρειάζεται, ανακατασκευάζουμε το υποδένδρο του.

- Μια απόσπαση στο IST γίνεται ως εξής: βρίσκουμε το x και το μετατρέπουμε ως εμβόλιμο και αντιστρέφουμε κατά 1 όλους τους κόμβους που βρίσκονται στο μονοπάτι απ' το x στη ρίζα. Έπειτα, βρίσκουμε τον υψηλότερο κόμβο του οποίου ο μετρητής έχει υπερβεί και ανακατασκευάζουμε το υποδένδρο του.

- Ένα IST με όρια a και b για ένα σύνολο $S = \{x_1 < x_2 < \dots < x_n\} \subseteq [a, b]$ με n στοιχεία αποτελείται από δύο πίνακες: τον $REP[1..n]$ με $n/2 \leq k \leq 2n/2$, ο οποίος περιέχει ένα δείγμα του file S (στα ιδανικά IST απαιτούμε το δείγμα να είναι ισοκαταμερισμένο) και τον $ID[1..n]$ με $m \leq k$, στον οποίο αποθηκεύουμε την αντίστροφη κατανομή των παιδιών (επίτευξη τελειών interrobate)

- Στην επιλογή του ευθετιού και ευθιού φαξίματος συμφωνούμε το στοιχείο x που ψάχνουμε με το $REP[j], REP[j+2^0], REP[j+2^1], REP[j+2^2], \dots$ μέχρι να βρεθεί κάποιο ώστε:

$REP[j+2^{k-1}] < x \leq REP[j+2^k]$. Με k το φαξίμο, ο χειρότερος χρόνος στο άσπασ γίνεται $O(\log n)$ και συνολικά στο IST λαμβάνουμε $O(\log^2 n)$.

- Για το φαξίμο, το IST παρουσιάζει ^{μικρή} χειρότερη συμπεριφορά αν τα βάρη κατά μήκος κάθε μονοπατιού απ' τη ρίζα στα φύλλα έχουν γεωμετρική πτώση, διότι έτσι προκύπτει λογαριθμικός βάθος. Για ενδέσσεις ή αποσβέσεις έχουμε καλή χειρότερη συμπεριφορά αν οι κόμβοι μεγέθους n μένουν αναδύονται για πολύ χρόνο, και οι δύο συμπεριφορές επιτυγχάνονται μέσω των αριθμών των κόμβων.

Έστω σύνολο $U = \{0, \dots, n-1\}$ και ένα hashtable (hash ταμπλέτα) $T[0, \dots, m-1]$. Τότε ορίζουμε μια hash συνάρτηση $h: U \rightarrow [0, \dots, m-1]$ τέτοια ώστε να αποθηκεύει ένα στοιχείο $x \in U$ στο $T[h(x)]$.

• Για παράδειγμα, έχουμε το $T = [0, 1, 2, 3, 4]$ ($m=5$) και δίδουμε μέσα σ' αυτό να αποθηκεύουμε το σύνολο $S = \{3, 15, 22, 24\}$ που είναι υποσύνολο του συνόλου $U = \{0, \dots, 99\}$. Χρησιμοποιούμε τη συνάρτηση $h(x) = x \bmod m \Rightarrow h(x) = x \bmod 5$.

0	15
1	
2	22
3	3
4	24

T

Είναι:

$$\begin{aligned} h(3) &= 3 \bmod 5 = 3 \\ h(15) &= 15 \bmod 5 = 0 \\ h(22) &= 22 \bmod 5 = 2 \\ h(24) &= 24 \bmod 5 = 4 \end{aligned}$$

Παρατηρούμε ότι δεν έχουμε καμία σύγκρουση (collision).
Αλλά, αν είχαμε να ενθέσουμε και το $x=20$, τότε $h(20) = 20 \bmod 5 = 0$. Στην θέση 0 όμως είναι αποθηκευμένο το 15. Άρα έχουμε πρόβλημα.

• Το παραπάνω πρόβλημα λύνεται με hashing with chaining (με αλυσίδες) και hashing with open addressing (ανοιχτές δεικνύσεις).

• Ανάλυση του hashing with chaining: Όλα τα κελιά του T είναι κεφαλές γραμμικών λιστών. Επομένως, η i-οστή λίστα θα περιέχει όλα εκείνα τα στοιχεία ενός συνόλου S για τα οποία $h(x) = i$. Για παράδειγμα, θεωρούμε τη hash ταμπλέτα $T = [0, 1, 2]$ ($m=3$) και το σύνολο $S = \{1, 3, 4, 7, 10, 17, 21\}$. Η hash συνάρτηση είναι η $h(x) = x \bmod m \Rightarrow h(x) = x \bmod 3$.

0	3 21
1	1 4 10
2	7 17

T

$$\begin{aligned} h(1) &= 1 \bmod 3 = 1 & h(7) &= 7 \bmod 3 = 1 & h(21) &= 21 \bmod 3 = 0 \\ h(3) &= 3 \bmod 3 = 0 & h(10) &= 10 \bmod 3 = 1 \\ h(4) &= 4 \bmod 3 = 1 & h(17) &= 17 \bmod 3 = 2 \end{aligned}$$

⊗ Στην μέση περίπτωση, η πολυπλοκότητα για ψάξιμο ενός στοιχείου $x \in S$, δηλαδή το να βρούμε την $h(x)$ και μετά να ψάξουμε το x στην αντίστοιχη λίστα $T[h(x)]$, απαιτεί σταθερό χρόνο, όσο και το μήκος της συγκεκριμένης λίστας. Το ίδιο συμβαίνει και για τις ελέγξεις και διαγραφές, αφού αυτές σε μια λίστα εκτελούνται σε $O(1)$ χρόνο.

Στη χειρότερη περίπτωση, η πολυπλοκότητα για ψάξιμο, ελέγξη και απόσβεση είναι $O(|S|)$. Αυτή εμφανίζεται όταν αποθηκεύει το σύνολο S αποθηκεύεται σε μία και μόνη λίστα.

• Μελετώντας μια αλυσίδα από η ενθέσεις, διαγραφές ή ψαξίματα, θα μάθουμε τις εξής παραδοχές:

- Η hash συνάρτηση h κατανέμει ισοφerrώς τα στοιχεία του συνόλου U πάνω στο διάστημα του $T[0, \dots, m-1]$, δηλαδή για όλα τα $i, j \in [0, \dots, m-1]$ ισχύει: $|h^{-1}(i)| = |h^{-1}(j)|$.
(αν π.χ. έχω σύνολο 9 στοιχείων και $m=3$, τότε 3 στοιχεία θα πάνε στο 0, 3 στο 1, 3 στο 2)
- Έχουμε ομοιόμορμη κατανομή στον είσοδο, δηλαδή όλα τα στοιχεία του U εμφανίζονται ισοferrώς σαν στοιχεία μιας πράξης, με πιθανότητα δηλαδή $1/|U|$.

⊗ Ορίζουμε ως $b = \frac{n}{m}$ τον παράγοντα κατανομής. Τότε, αν πάρουμε το μέσο πλήθος των στοιχείων της λίστας i θα πρέπει $m \cdot (\text{μέσο πλήθος}) = n$, και το $b = \frac{n}{m}$ θα είναι αυτό το μέσο πλήθος.

- Πολυπλοκότητα μέσης περίπτωσης για n πράξεις είναι $O(1 + \frac{b}{m}|n|)$, όπου $b = \frac{n}{m}$ είναι ο μέσος παράγοντας κατανομής. Το κατανομημένο κόστος (μέσο) είναι τότε $O(1)$.
- Το μέσο μήκος της μεγαλύτερης αλυσίδας είναι $O(\frac{\log n}{\log b})$ όταν το $b = \frac{n}{m} \leq 1$.

• Ανάλυση του hashing with open addressing: Εδώ όλα τα στοιχεία του συνόλου U αποθηκεύονται στο T. Δηλαδή, όλα τα κελιά του T θα περιέχουν ένα στοιχείο ή θα είναι κενά. Πότε τέτοιο γίνεται εμείς με το να θεωρούμε ότι κάθε στοιχείο $x \in U$ ορίζει μια αλυσίδα από hash συνάρτησεις $h(0, x), h(1, x), \dots$ οι οποίες αποτελούν θέσεις στο T. Με άλλα λόγια, ενθέτουμε ένα στοιχείο x αποθηκεύοντας αυτό την αλυσίδα μέχρι να βρεθεί μια ελεύθερη θέση για το x. Αριθμείς με τον ίδιο τρόπο λειτουργεί και το ψάξιμο.

- Έχουμε τα hash συνδυασμών $h(i, x) = [h_1(x) + i \cdot h_2(x)] \bmod m$, όπου $h_1(x) = x \bmod 7$, $h_2(x) = 1 + x \bmod 4$.
Διαθέτουμε για hash ταμπλέτα $T = [4, 2, 3, 4, 5, 6]$ ($m = 7$) και ακολουθώντας αυτήν την ακολουθία, θέ-
λουμε να αποθηκεύσουμε στο T το σύνολο $S = \{2, 17, 6, 9, 15, 13, 10\}$. Έχουμε:

0	
1	15
2	9
3	3
4	
5	17
6	6
	T

$$h(0,3) = h_1(3) \bmod 7 = [3 \bmod 7] \bmod 7 = 3$$

$h(0, 17) = h_1(17) \bmod 7 = [17 \bmod 7] \bmod 7 = 3 \bmod 7 = 3 \Rightarrow$ Existenz einer Funktion, die aus den bits für ein aus den bits für $i=1$.

$$h(4, 17) = [h_1(17) + h_2(17)] \bmod 7 = [17 \bmod 7 + 1 + 17 \bmod 4] \bmod 7 = (3 + 1 + 1) \bmod 7 = 5$$

$$h(0,6) = h_1(6) \bmod 7 = (6 \bmod 7) \bmod 7 = 6$$

$$h(0,9) = h_1(9) \bmod 7 = 2$$

$$h(0, 13) = (13 \bmod 7) \bmod 7 = 6 \text{ (collision)}$$

$$h(0, 15) = h_1(15) \bmod 7 = (1)$$

$$h(1, 13) = [h_1(13) + h_2(13)] \bmod 7 = \dots \text{collision} \dots$$

$$h(\underline{4}, 13) = [13 \bmod 7 + 4 + \cancel{4} \cdot 13 \bmod 4] \bmod 7 =$$

$$= (6+4+52 \bmod 4) \bmod 7 = \text{collision..}$$

$h(0, 10) = \dots = \text{collision} \dots$

- Δεν χρησιμοποιείται τον επιπλέον χώρο για αποθήκευση δεδομένων. Στην διαφραγή όμως έχουμε προβλημα με το απειδελας οβιότο ενός στοιχείου, γιατί τότε χαλάει η δομή και δε γίνεται σωστά η ένθεση και το ψάξιμο. Γι' αυτό ένα κελί θεωρούμε ότι έχει 3 καταστάσεις: ελεύθερο, καταληφθέν ή οβιότο. Δηλαδή, στη διαφραγή καταγράφουμε ένα στοιχείο ως οβιότο (lazy deletion).

- Αποκρίνονται με τις επλήρωσεις για εύθετο $= 0 \left(\frac{1}{1-\beta} \right)$ όπου $\beta < 1 \Rightarrow \frac{n}{m} < 1 \Rightarrow n < m$.

πασιφο = $O\left(\frac{1}{8} \ln \frac{1}{1-\delta}\right)$, όπου $n \leq m$

Στη χειρότερη περίπτωση και στις δύο περιπτώσεις έχουμε πολυπλοκότητα $O(n)$

- Συνίθως, το open addressing απαιτεί να επιλέξουμε το m μεγαλύτερο d ήως μπορεί να αντιστοιχιστεί μια ταύτιση χάρου. Γι' αυτό μια συνήθης ταυτίση είναι να επιλέξουμε αρχικά μικρό m , δηλαδή δι' μικρό μέγεθος του T και να το διπλασιάζουμε ϕ ταν το n γίνει αρκετά μεγάλο.

Για να αναπαραστήσουμε μια τέτοια διαδρομή για ακολουθία από hash-ταβλιέρες T_0, T_1, T_2, \dots με αντίστοιχες hash αναδρομικές h_0, h_1, h_2, \dots όπου $|T_0| = m$ και $|T_i| = 2^i \cdot m$. Οπότε, για τον παραπάνω υπολογισμό R_i έχουμε τα εξής:

$T_i = \begin{cases} T_{i-1} & , \quad b_i = 1 \text{ (Α)} \\ T_{i-1} & , \quad b_i = \frac{1}{u} \text{ (Β)} \end{cases}$ Όταν ο b γίνει ίσος με 1, δηλαδή $n=m$ και ο T_i έχει j φορές με 2^i στοιχεία τότε συμπληρώθηκε των $T_{i+1} = 2^{i+1}$ και ενδεχομένως να προταφέντα στοιχεία.

- Τα έξοδα μετάβασης για να πάμε απ' τον Τι \leadsto Τιτι είναι $O(2^i \cdot m)$. Τότε προκύπτει ότι $\text{bit}_i = \frac{1}{2}$ και πρέπει να αυτεξοδωσούμε τουλάχιστον $\frac{1}{2} \cdot (2^i \cdot m) = \frac{1}{4} \cdot (2^{i+1} \cdot m)$ πράξεις στον Τιτι μέχρι να γίνει ο παράγοντας των $\text{bit}_i = \frac{1}{4}$ και να παραχθούν νέα έξοδα μετάβασης. Με κατανομή αυτών των εξόδων έχουμε καταμετρημένο κόστος $\frac{2^i \cdot m}{2^{i+1} \cdot m} = 2 = O(4)$

- Τα έξοδα μετάβασης για να πάμε από τον $T_i \leadsto T_{i-1}$ είναι $O(\frac{1}{4} \cdot 2^{i-m})$. Τότε προκύπτει ότι $\beta_{i-1} = \frac{1}{2}$ και πρέπει να ευντελέσουμε πυλώνισαν $\frac{1}{8} 2^{i-m} = \frac{1}{4} \cdot (2^{i-m})$ πράξας στο T_{i-1} μέχρι να γίνει απαράφορος του $\beta_{i-1} = \frac{1}{4}$. Ανάλογα έχουμε κατασκευασμένο κόσος $O(1) = 2$.


• Extendible hashing: χρησιμοποιείται όταν η κύρια μνήμη δεν είναι αρκετά μεγάλη για να αποθηκεύει δένδρα το hash πίνακα. Γι' αυτό μεταφέρουμε το hashing σε βοηθητική μνήμη (secondary memory).

Θεωρούμε ότι η βοηθητική μνήμη είναι χωρισμένη σε καύδους (buckets) ορισμένου μεγέθους A στους οποίους αποθηκεύουμε δεδομένα μεγέθους B . Αν είναι $S \subseteq U$ το υποσύνολο του σιμπόλου U που θέλουμε να αποθηκεύσουμε και αν είναι $d(S)$ το βάθος του S σε σχέση με το B , όπου B είναι το μέγεθος ενός bucket, τότε αποθηκεύουμε όλα τα στοιχεία σε buckets τα οποία αποτελούν φύλλα ενός δένδρου T με βάθος $d(S)$. Τότε, το ψάξιμο στα φύλλα επιτυγχάνεται σε $O(1)$ χρόνο. (το μέσο πλήθος για n buckets είναι $O(\frac{\log n}{\log A})$)

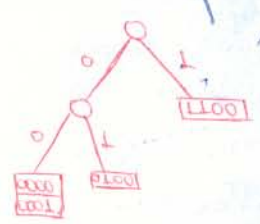
• Το extendible hashing χρησιμοποιεί μια hash ταμπλέτα $T = [0, \dots, 2^{d(S)} - 1]$ που ονομάζεται directory και μερικά buckets για να αποθηκεύεται τα στοιχεία του συνόλου S .

Για παράδειγμα, θεωρούμε buckets μεγέθους $b=2$, μια συνάρτηση $h: U \rightarrow \{0,1\}^k$, $k \in \mathbb{N}$ και το σύνολο S του οποίου θα αποθηκεύσουμε τα στοιχεία, και $h(S) = \{0000, 0001, 0100, 1100\}$.

Φαίνεται το αντίστοιχο tree:



Βλέπουμε ότι ένα bucket σε χρησιμοποιείται καλύπτει, οπότε θα συμπιέσουμε το tree (compression).



Εδώ δαμε τον κλάδο χωρίς παιδί και έπειτα τον αντίστοιχο κώδικα.

Παρατηρούμε ότι το ύψος του tree είναι 2 άρα $d(S)=2$. Οπότε, σχεδιάζουμε τη δέση:

00	→	0000 0001 (2)
01	→	0100 (2)
10	→	
11	→	1100 (1)
T		

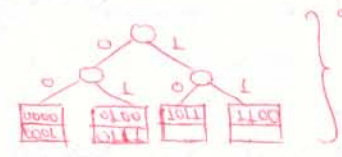
Τα αντίστοιχα βάθη των buckets, που φαίνονται στις παρενθέσεις, είναι τα τμήματά βάθη των buckets στο tree που έγινε compressed. Γι' αυτό θα ισχύει, ότι αν είναι $2^{d(S)-\tau}$ ο αριθμός των στοιχείων του T που δείχνουν σε ένα bucket, τότε το τ είναι το τμήμα βάθος του bucket.

(π.χ. για το τελευταίο bucket είναι: $2^{d(S)-\tau} = 2 \Rightarrow d(S) - \tau = 1 \Rightarrow 2 - \tau = 1 \Rightarrow \tau = 1$)

• Για να ενδεύουμε ένα στοιχείο κοιτάμε αν αυτό χωράει σε ένα απ' τα υπάρχοντα buckets. Αν δε χωράει τότε διπλασιάζουμε το μέγεθος του T και κατασκευάζουμε το tree απ' την αρχή.

→ Στο προηγούμενο παράδειγμα, αν θέλουμε να ενδεύουμε το στοιχείο 0111 βλέπουμε ότι το συγκεκριμένο bucket 0100 έχει χώρο, οπότε το ενδεύουμε χωρίς πρόβλημα.

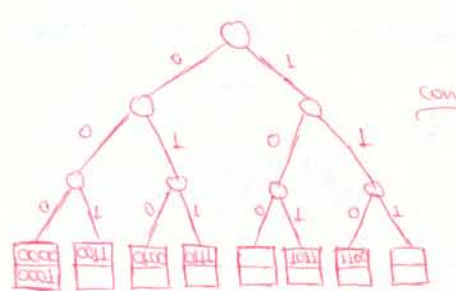
→ Αν, τώρα, θέλουμε να ενδεύουμε το στοιχείο 1011 θα πρέπει να κάνουμε το tree όπως ήταν πριν, δηλαδή decompressed:



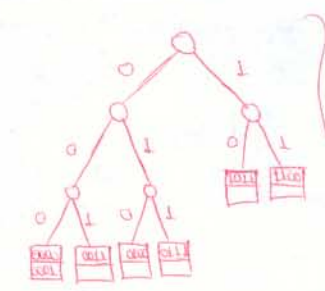
οπότε ο T γίνεται:

00	→	0000 0001 (2)
01	→	0100 0111 (2)
10	→	1000 1001 (2)
11	→	1100 (2)
T		

→ Αν, όμως, θέλουμε να ενδεύουμε το στοιχείο 0011 (δεν χωράει) θα πρέπει να κάνουμε νέο tree και να διπλασιάσουμε το T :



compression



000	→	0000 0001 (3)
001	→	0011 (3)
010	→	0100 (3)
011	→	0111 (3)
100	→	1000 1001 (2)
101	→	1010 1011 (2)
110	→	1100 (2)
111	→	
T		

• Για να διαγράψουμε ένα στοιχείο, κοιτάμε αν αρχικά ο κώδικας είναι μεγαλύτερος. Αν είναι, αφαιρούμε το στοιχείο (διαγράφουμε). Αν δεν είναι μεγαλύτερος και έχει 1 ή μόνο στοιχείο τότε μειώνουμε το βάθος του δένδρου στο μισό και κάνουμε τις απαραίτητες διευθετήσεις στον T (π.χ. θα δείχνουν δηλαδή τα στοιχεία μετά τη μείωση)

UNION-FIND

Θεωρούμε ένα σύνολο $U = \{0, 1, \dots, N-1\}$ και το διαχωρίζουμε σε υποσύνολα. Ένα μεταξύ τους σύνολο δύ που να μην έχουν κοινά στοιχεία το ένα με το άλλο. Τότε ορίζουμε τις πράξεις $\text{find}(x)$, που επιστρέφει το υποσύνολο εκείνο στο οποίο βρίσκεται το x και $\text{union}(A, B, C)$ η οποία επιστρέφει ένα σύνολο $G \subseteq U$ το οποίο είναι συνένωση των A και B .

- Χρησιμοποιούμε έναν πίνακα $\text{IS-in}[0, \dots, N-1]$, όπου $\text{IS-in}[i]$ είναι το όνομα του υποσυνόλου που περιέχει το i . Τότε η $\text{find}(x) = \text{IS-in}[x]$ εκτελείται σε $O(1)$ χρόνο (constant) αλλά η $\text{union}(A, B, C)$ σε χρόνο $O(N)$, αφού θα πρέπει να διασχίσουμε όλο τον πίνακα και να αντιστοιχίσουμε τα A, B .

Και τέτοιο μπορεί να αποφευχθεί αν κάθε σύνολο διατηρεί μια λίστα των στοιχείων που περιέχει. Τότε για $\text{union}(A, B, C)$ θα πρέπει να διατρέξουμε μόνο τις λίστες των A και B και όχι όλο τον πίνακα. Σ' αυτή την περίπτωση έχουμε χρόνο $O(|A| + |B|)$. Ο χρόνος αυτός βελτιώνεται σε $O(\min\{|A|, |B|\})$ με την εφαρμογή του weighted union rule, κατά το οποίο σε κάθε λίστα ενός υποσυνόλου αποθηκεύω και το μέγεθός της, οπότε διαλέγουμε μόνο τη μικρότερη λίστα και την ενώνουμε στην άλλη.

Στην προηγούμενη λίστα υπάρχει ένα πρόβλημα. Αν οι συνένωση των υποσυνόλων να ονομάζονται G για $\text{union}(A, B, G)$ αυτή ονομάζεται A (αν π.χ. $|A| > |B|$).

- Σ' αυτή την περίπτωση χρησιμοποιούμε 2 απεικονίσεις πινάκων, τους MAPIN και MAPOUT . Θα είναι

MAPIN : εξωτερικό όνομα \rightarrow εσωτερικό όνομα
 MAPOUT : εσωτερικό όνομα \rightarrow εξωτερικό όνομα

Συλλογή για το παραπάνω θα είναι $\text{MAPOUT}[A] = C$

Εξωτερικό όνομα υποσυνόλου με εσωτερικό όνομα A .

Οπότε, η find θα δίνεται ως $\text{find}(x) = \text{MAPOUT}[\text{IS-in}[x]]$

\rightarrow εσωτερικό όνομα του υποσυνόλου που περιέχει το x

- Για την πράξη $\text{union}(A, B, C)$, αρχικά ελέγχουμε το μέγεθος των λιστών A και B (χρόνος $O(1)$), έπειτα αλλάζουμε το εσωτερικό όνομα του μικρότερου υποσυνόλου και συνδέουμε τις λίστες (χρόνος $O(1)$) και τελικά ενημερώνουμε τα MAPOUT , MAPIN και SIZE (χρόνος $O(1)$). Ο $\text{SIZE}[A]$ δείχνει το μέγεθος του υποσυνόλου με εξωτερικό όνομα A .

- Μια ακολουθία από m Finds και $(n-1)$ Unions έχει συνολικό κόστος $O(m + n \log n)$.

- Ένα find κοστίζει $O(1)$ γ'αυτό έχουμε $O(m)$ για m Finds.
- Ένα $\text{union}(A, B, C)$ κοστίζει $O(\min\{|A|, |B|\})$ ή $O(\min\{\text{size}[A], \text{size}[B]\}) \leq C \cdot \min\{\text{size}[A], \text{size}[B]\}$, για μια σταθερά C .

Για $n-1$ Unions θα ισχύει ότι $|A| \leq n$ και το σύνολο C θα έχει τονλάχιστον διπλάσιο μέγεθος απ' τα A και B . Γροφένως: $2 \cdot 2 \cdot 2 \dots 2 \leq n \Rightarrow 2^i \leq n \Rightarrow i \leq \log n$, και αν το στοιχείο $x \in A$ τότε αυτό θα αλλάζει μέρος σε ανατάξεις κατά $n/2^i$ φορές. Στη χειρότερη περίπτωση θα είναι $2^i = n$ και $n/2^i = 1 \Rightarrow i = \log n$. Δηλαδή ένα Union κοστίζει $O(C \log n)$. Για $n-1$ Unions μπορεί να λάβουν μέρος το πολύ $2 \cdot (n-1)$ στοιχεία. Άρα: $2(n-1) \cdot C \log n \Rightarrow (n-1)$ Unions κοστίζουν στη χειρότερη περίπτωση $O(n \log n)$. Έτσι το κατωφλιωμένο κόστος για ένα Union θα είναι $O(\log n)$.

- Το union-find πρόβλημα μπορεί να λυθεί και διαφορετικά. Μπορούμε να παραστήσουμε κάθε υποσύνολο σαν ένα δένδρο και ως όνομα στο κάθε υποσύνολο να χρησιμοποιούμε τη ρίζα του. (υλοποίηση με 3 πίνακες: Father , $\text{Name}[0, \dots, N-1]$, $\text{Root}[j]$).

Σ' αυτή την περίπτωση το union κοστίζει $O(1)$ και το find $O(n)$ χρόνο. Μια πρώτη βελτίωση είναι να εφαρμόσουμε weighted union rule έτσι ώστε οι ρίζες των δένδρων με λιγότερους κόμβους να δείχνουν στις ρίζες των δένδρων με τους περισσότερους (δηλαδή να γίνονται παιδιά τους).

π.χ.: $\text{Union}(1, 2, A)$ $\text{Union}(3, 4, B)$ $\text{Union}(A, B, C)$ $\text{Union}(5, 6, D)$ $\text{Union}(7, 8, E)$ $\text{Union}(A, E, F)$ $\text{Union}(C, F, G)$



- Αυτά των περιπτώσεων για ακολουθία από $n-1$ Unions και m finds υποτίθεται $O(n + m \log n)$
- Αν είναι F το δέντρο που παίρνουμε μετά από $n-1$ Unions τότε ορίζουμε ως $\text{rank}(x)$ το ύψος του x μέσα στο F και $\text{weight}(x)$ τον αριθμό των κόμβων στο υποδένδρο που ελέγχει το x .
- Θέλουμε να δείξουμε ότι $\text{weight}(x) \geq 2^{\text{rank}(x)}$ και $\text{rank}(x) \leq \log n$
- $\text{weight}(x) \geq 2^{\text{rank}(x)}$: Αποδεικνύουμε επαγωγικά ως εξής: $\text{rank}(x)=0 \Rightarrow 2^{\text{rank}(x)}=1 \leq \text{weight}(x)$.
Για $\text{rank}(x) > 0$, αν είναι y ο γιος του x με $\text{rank}(y) = \text{rank}(x) - 1$, τότε όταν ο y έγινε γιος του x τότε ο x είχε τουλάχιστον δύο απογόνους όπως είχε και ο y (weighted union rule).



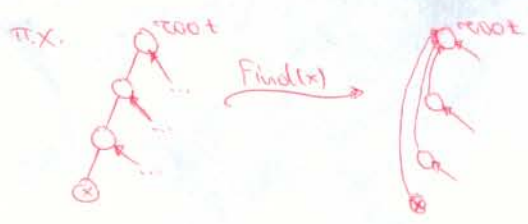
21

Οπότε $\text{weight}(x) \geq \text{weight}(x') + \text{weight}(y) \Rightarrow (\text{weight}(y) \leq \text{weight}(x)) \Rightarrow \text{weight}(x) \geq 2 \text{weight}(y) \Rightarrow \text{weight}(x) \geq 2 \cdot 2^{\text{rank}(y)}$
 $\Rightarrow \text{weight}(x) \geq 2 \cdot 2^{\text{rank}(x)-1}$ (επαγωγική υπόθεση)

Τότε: $\text{weight}(x) \geq 2^{\text{rank}(x)+1} \Rightarrow \boxed{\text{weight}(x) \geq 2^{\text{rank}(x)}}$

- $\text{rank}(x) \leq \log n$: Με $n-1$ Unions μπορεί να παραχθούν σίγουρα μέγιστος το ποσό $\log n$ ($\leq n$), οπότε $\text{weight}(x) \leq n \Rightarrow 2^{\text{rank}(x)} \leq \text{weight}(x) \leq n \Rightarrow \boxed{\text{rank}(x) \leq \log n}$
 ή $\text{weight}(\text{root}) \geq 2^{\text{rank}(\text{root})} \Rightarrow n \geq 2^{\text{rank}(\text{root})} \Rightarrow \text{rank}(\text{root}) \leq \log n$

- Εξαιτίας του $\text{rank}(x) \leq \log n$ έκαστος find υποτίθεται το ποσό $\log n$ φορές. Ένα Union στοιχίζει $O(1)$
- Μια δεύτερη βελτίωση είναι το path compression: κατά την εκτέλεση ενός find(x) στο μονοπάτι από το x ως προς την ρίζα παίρνουμε τους δεικτες και τους βάσουμε να δείχνουν στη ρίζα. Κατά τέτοιο τρόπο τα επόμενα finds γίνονται γρηγορότερα.



- Αν στο Union-Find πρόβλημα εφαρμόσουμε weighted union rule και path compression μαζί τότε το συνολικό κόστος για m ακολουθίες από $n-1$ Unions και m finds είναι $O(m \cdot \alpha(m, n))$ όπου $\alpha(m, n) = \min \{ z \geq 1; A(z, 4 \log m) > \log n \}$ είναι μια παραλλαγή της συνάρτησης Ackermann.