

Δομές Δεδομένων: Εργασία 2

ΑΠΟ ΒΑΣΙΛΗΣ ΠΑΠΑΔΗΜΑΣ (3220150) & ΜΑΡΙΟΣ ΜΑΤΣΑ (3220120)

Μέρος Α: Influenza_k.java

Στο αρχείο Influenza_k.java έχουμε υλοποιήσει την συνάρτηση quicksort(City[] citys, int start, int end) η οποία αποτελεί μια εξειδικευμένη εκδοχή του αλγορίθμου QuickSort, προκειμένου να μπορούμε να επεξεργαστούμε τα δεδομένα τύπου City.

Χαρακτηριστικά αλγορίθμου:

- 1) Time complexity: $O(n \log n)$ - Space complexity: $O(n)$
- 2) Τελευταίο στοιχείο ως pivot
- 3) Ταξινόμηση τέτοια ώστε το πρώτο στοιχείο να έχει το μικρότερο πλήθος κρουσμάτων ανά 50,000 κατοίκους (Σε περίπτωση ισοβαθμίας μεταξύ πόλεων, θεωρούμε πιο υψηλά στην κατάταξη το στοιχείο που προηγείται αλφαβητικά και σε περίπτωση που και τα ονόματα ταυτίζονται, προηγείται αυτό με το μικρότερο ID).

Όταν τρέχουμε το αρχείο Influenza_k.java από το terminal θα πρέπει να δώσουμε ως όρισμα το k και το file path που περιέχει τις πόλεις και τα δεδομένα που θα χρησιμοποιήσουμε. Σε περίπτωση που δοθούν περισσότερα ή λιγότερα ορίσματα το πρόγραμμα τερματίζει με exit code 1 και ένα αντίστοιχο μήνυμα εμφανίζεται στον χρήστη. Αν το k δεν είναι integer το πρόγραμμα και πάλι τερματίζει, ενώ αν όλα τα ορίσματα είναι σωστά τότε τα αποθηκεύουμε απλώς σε 2 μεταβλητές και τα χρησιμοποιούμε κανονικά στο πρόγραμμα.

Μέρος Β: ΑΤΔ ουράς προτεραιότητας

Η συνάρτηση remove(int id) έχει Time complexity $\rightarrow O(\log n)$. Προκειμένου όμως να είναι κάτι τέτοιο δυνατό και δεδομένου ότι η συνάρτηση sink(int index) και swim(int index) που χρησιμοποιούμε έχουν επίσης time complexity $\rightarrow O(\log n)$, τότε η αναζήτηση του στοιχείου προς διαγραφή θα πρέπει να γίνεται σε χρόνο καλύτερο του n. Στην προκειμένη περίπτωση χρησιμοποιούμε έναν πίνακα με 1000 θέσεις για τον οποίο ισχύει το εξής:

- 1) Key \rightarrow Το id της χώρας (π.χ μια χώρα με id = 20 θα βρίσκεται στην θέση 20 του πίνακα)
- 2) Value \rightarrow Η θέση του στοιχείου στο οποίο αντιστοιχεί το id στον πίνακα heap (π.χ αν η χώρα με id 20 βρίσκεται στην θέση 3 του πίνακα heap, τότε θα ισχύει id[20]=3)

Με αυτό τον τρόπο βρίσκουμε το στοιχείο που ψάχνουμε σε χρόνο $O(1)$, το αντικαθιστάμε με το τελευταίο και εκτελούμε sink και swim (δεν ξέρουμε αν το νέο στοιχείο θα είναι μικρότερο του πατέρα/μεγαλύτερο από τα παιδιά του) για να ξαναδημιουργήσουμε ένα σωστό heap.

Φυσικά η επιλογή των 1000 θέσεων δεν είναι αυθαίρετη, καθώς τα idE[1, 999] και άρα μπορεί να έχουμε το πολύ 999 ids (η επιπλέον θέση 0 δεν χρησιμοποιείται γιατί δεν υπάρχει id = 0).

Μέρος Γ: DynamicInfluenza_k_withPQ.java

Στο DynamicInfluenza_k_withPQ.java, αρχικά μετράμε το μέγεθος του αρχείου εισόδου και μετά ζητάμε από τον χρήστη να εισάγει το k. Επειτα ελέγχουμε αν $k > N$, όπου N ο αριθμός γραμμών του αρχείου εισόδου. Μετά, δημιουργούμε ένα PQ cities με μέγεθος 2k και Comparator της τάξης NegativeComparator (η οποία ουσιαστικά ταξινομεί τα στοιχεία με την αντίστροφη σειρά από την κανονική, δηλαδή το min στοιχείο θα είναι το max στοιχείο κ.ο.κ). Τότε, κάνουμε ένα loop που διαβάζει όλες τις γραμμές του αρχείου εισαγωγής και δημιουργεί για κάθε γραμμή ένα αντικείμενο City με τα αντίστοιχα δεδομένα. Εάν η PQ έχει λιγότερα από k στοιχεία τότε το εισάγει σε αυτή. Αλλιώς, το εισάγει μόνο εάν έχει μεγαλύτερη προτεραιότητα από το στοιχείο της PQ με την μικρότερη προτεραιότητα, το οποίο σύμφωνα με τον NegativeComparator θα είναι το PQ.min(). Σε αυτή την περίπτωση επίσης αφαιρεί το min ώστε η PQ να έχει k αντικείμενα.

Αφού τελειώσει το loop, σε ένα άλλο loop εξάγουμε τα k στοιχεία της PQ (κάθε φορά με το `getmin()`) και εισάγουμε το όνομα τους σε έναν `String[] leaderboard` με την αντίστροφη σειρά. Έτσι, όταν στο επόμενο και τελευταίο loop εκτυπώνουμε τα περιεχόμενα του leaderboard, τα ονόματα των πόλεων είναι σε αύξουσα σειρά πυκνότητας χρουσμάτων.

Επομένως, για την ανάλυση της πολυπλοκότητας το ουσιαστικό κομμάτι του κώδικα που πρέπει να αναλύσουμε είναι το loop που επεξεργάζεται το input και εισάγει ή μη τα αντικείμενα City στην PQ. Σε αυτό το loop έχουμε N επαναλήψεις. Κάθε επανάληψη έχει σταθερό κόστος το splitting του String με regex, την δημιουργία του αντικειμένου City και την ρύθμιση των πεδίων του. Όλα αυτά είναι $O(1)$. Για τις πρώτες k επαναλήψεις κάνουμε πάντα insert το αντικείμενο στη PQ, αυτό είναι $O(\log n)$. Στις επόμενες $N-k$, βρίσκουμε το min αντικείμενο ($O(1)$ λόγω της υλοποίησης μας) και εάν έχει μικρότερη προτεραιότητα από το τρέχον city τότε το αφαιρούμε $O(\log n)$ και εισάγουμε το τρέχον $O(\log n)$. Άρα το συνολικό loop έχει worst-case complexity $O(N \log n)$. Στην πραγματικότητα, είναι $O(N)$ μετά από την επανάληψη που θα εισαχτεί το στοιχείο που έχει globally την k -οστή υψηλότερη προτεραιότητα. Επομένως, άμα αυτό συμβεί νωρίς (περιπτώσεις όπου το k είναι αρκετά μικρότερο του συνολικού αριθμού πόλεων) τότε η πολυπλοκότητα θα πλησιάζει $O(N)$ (best-case άμα τα πρώτα k στοιχεία του input είναι οι k πόλεις με την μικρότερη πυκνότητα).

Όταν τρέχουμε το αρχείο `DynamicInfluenza_k.java` από το terminal θα πρέπει να δώσουμε ως όρισμα το k και το file path που περιέχει τις πόλεις και τα δεδομένα που θα χρησιμοποιήσουμε. Σε περίπτωση που δοθούν περισσότερα ή λιγότερα όριασμα το πρόγραμμα τερματίζει με exit code 1 και ένα αντίστοιχο μήνυμα εμφανίζεται στον χρήστη. Αν το k δεν είναι integer το πρόγραμμα και πάλι τερματίζει, ενώ αν όλα τα όριασμα είναι σωστά τότε τα αποθηκεύουμε απλώς σε 2 μεταβλητές και τα χρησιμοποιούμε κανονικά στο πρόγραμμα.

Μέρος Δ: Dynamic_Median.java

Στο `Dynamic_Median.java`, δημιουργούμε δυο PQ (hi και lo) με μέγεθος 500 (ώστε να μην γίνει ποτέ `resize`). Η lo χρησιμοποιεί τον `NegativeComparator` και η hi τον `PositiveComparator`. Μετά κάνουμε ένα loop, διαβάζοντας σε κάθε επανάληψη μια γραμμή του αρχείου εισόδου και δημιουργώντας ένα αντικείμενο City με τα αντίστοιχα στοιχεία. Τότε, αυτό που θέλουμε να κάνουμε είναι να εισάγουμε το αντικείμενο στην hi ή την lo , με στόχο όλες οι k πόλεις που έχουμε διαβάσει στην k -οστή επανάληψη να είναι διαμερισμένες στις hi και lo έτσι ώστε όλες οι πόλεις στην hi να έχουν υψηλότερη προτεραιότητα από όλες τις πόλεις στην lo , και επίσης αν ο k είναι άρτιος να ισχύει $size(lo) = size(hi) = k/2$ (συνθήκη 1) και αν ο k είναι περιττός να ισχύει $size(lo) = size(hi) + 1 = \lceil k/2 \rceil$ (συνθήκη 2). Ωστε να επιτευχτεί αυτό:

- αν και οι δύο PQ είναι κενές (πρώτη επανάληψη), το εισάγουμε στην lo . Ισχύει η συνθήκη 2.
- αν η hi είναι κενή και η lo περιέχει ένα στοιχείο (δεύτερη επανάληψη), εισάγουμε το τρέχον στην hi άμα έχει ψηλότερη προτεραιότητα από αυτό στην lo , αλλιώς τους αλλάζουμε σειρά. Ισχύει η συνθήκη 1.
- στις υπόλοιπες επαναλήψεις, εισάγουμε το στοιχείο στην lo άμα έχει χαμηλότερη προτεραιότητα από το στοιχείο στην lo με την υψηλότερη προτεραιότητα, αλλιώς το εισάγουμε στην hi . Τότε, αν το k είναι άρτιος ($(hi.size() + lo.size()) \% 2 == 0$), στην προηγούμενη επανάληψη ήταν περιττός, άρα αρχικά ίσχυε η συνθήκη 2, και η συνθήκη 1 δεν θα ισχύει μόνο άμα προσθέσαμε στο hi (δηλαδή $lo.size() == hi.size() + 2$, οπότε αφαιρούμε από την lo το στοιχείο με την υψηλότερη προτεραιότητα (λόγω `NegativeComparator`) και το προσθέτουμε στην hi για να έχουν ίδιο μέγεθος. Αλλιώς, άμα το k είναι περιττός τότε στην προηγούμενη επανάληψη ήταν άρτιος άρα ίσχυε η συνθήκη 1, και η συνθήκη 2 δεν θα ισχύει τώρα μόνο εάν προσθέσαμε στο hi ($lo.size() < hi.size()$) άρα αφαιρούμε το στοιχείο του hi με την χαμηλότερη προτεραιότητα και το προσθέτουμε στην lo ώστε να ισχύει η συνθήκη 2.

Πρακτικά λοιπόν, κάθε επανάληψη έχει κόστος $O(\log n)$ και συνολικά το πρόγραμμα μας έχει πολυπλοκότητα $O(N \log n)$, αφού ο ίδιος ο υπολογισμός του mean είναι ($O(1)$): όταν θέλουμε να υπολογίσουμε το mean, στο τέλος οποιασδήποτε επανάληψης, ξέρουμε ότι θα είναι το `lo.min()`.

Όταν τρέχουμε το αρχείο `Dynamic_Median.java` από το terminal θα πρέπει να δώσουμε ως όρισμα το file path που περιέχει τις πόλεις και τα δεδομένα που θα χρησιμοποιήσουμε. Σε περίπτωση που δοθούν περισσότερα ή λιγότερα ορίσματα το πρόγραμμα τερματίζει με exit code 1 και ένα αντίστοιχο μήνυμα εμφανίζεται στον χρήστη.