

ΜΥΕ041 - ΠΛΕ081: Διαχείριση Σύνθετων Δεδομένων

ΕΡΓΑΣΙΑ 1 – Αποτίμηση ερωτημάτων

Νεφέλη-Ελένη Κατσυλέρου

ΑΜ : 4385

Μέρος 1 (group-by with aggregation)

Αρχικά απο τη γραμμή εντολών παίρνουμε το `input_file` του `csv` αρχείου που θέλουμε ως είσοδο καθώς επίσης τα `attributes` για την ταξινόμηση και τη συνάθροιση και τη συνάρτηση `aggregation_function` (`max`, `min`, `sum`). Μετά καλώ τη συνάρτηση `group_by_aggregation`.

Η συνάρτηση `group_by_aggregation(input_file, clustering_attribute, aggregation_attribute, aggregation_function)` διαβάζει δεδομένα από το αρχείο `csv` που δίνεται μέσω του `input_file` και τα μετατρέπει σε μια λίστα απο `ints`. Διαβάζει κάθε γραμμή του αρχείου `csv` , τη μετατρέπει σε μια λίστα ακεραίων και την βάζει στη λίστα `data`. Στη συνέχεια καλείται η συνάρτηση `my_merge_sort` με τα `data` αυτά και τα υπόλοιπα ορίσματα που απαιτούνται.

Η συνάρτηση `my_merge_sort(data, clustering_attribute, aggregation_attribute, aggregation_function)` χρησιμοποιείται για να χωρίζει το σύνολο απο τα `data` που δίνονται ως `input` σε 2 υποσυνολα `left` και `right`. Μετα καλει αναδρομικα τον εαυτο της για καθε υποσυνολο μεχρι τα δεδομένα να έχουν μόνο ένα στοιχείο ή είναι κενά. Τότε καλει την `my_merge` για να κανει απο “κάτω” προς τα “πάνω” `merge` τα `left` και `right` υποσυνολα.

Στη συνάρτηση `my_merge(left, right, clustering_attribute, aggregation_attribute, aggregation_function)` έχω αρχικοποιήσει το `i` που είναι ο `counter` για το `left` και το `j` που είναι ο `counter` για το `j`, το `result=[]` είναι μια άδεια λίστα όπου θα αποθηκεύσω το `merge` αποτέλεσμα και το `aggregation_dictionary = {}` είναι ένα λεξικό για να αποθηκεύσω `aggregation values` που αντιστοιχούν σε κάθε `clustering attribute`. Τα κλειδιά του λεξικού αντιπροσωπεύουν μοναδικά `clustering attribute values`. Οι τιμές του λεξικού είναι λίστες που περιέχουν τις αντίστοιχες τιμές των `aggregation attribute`. Οι τιμές αυτές συμπληρώνονται με βάση τις τιμές των `clustering attribute` που γίνονται `count`.

Στη συνέχεια επαναληπτικά περνάει και απο τα δυό μισά και συγκρίνει τα `clustering attributes` για να κάνει κατάλληλα `merge`. Ανακτά τις τιμές των `clustering attributes` απο το αριστερό και το δεξί μισό για σύγκριση και ανάλογα με τη σύγκριση αποφασίζει ποιά στοιχεία θα γίνουν `merge`. Αφού ολοκληρωθεί η διαδικασία του `merge` , οι τιμές συνάθροισης υπολογίζονται απο τις λίστες του λεξικού και εγγράφονται στο ζητούμενο αρχείο εξόδου (`O1.csv`).

Μέρος 2(merge join)

Αρχικά ορίζω τη συνάρτηση `my_merge_join` που λαμβάνει τρεις παραμέτρους: `R_file`, `S_file` και `output_file`. Αυτές οι παράμετροι αντιπροσωπεύουν τις διαδρομές προς τα δύο αρχεία CSV εισόδου `R_file` και `S_file` και το αρχείο CSV εξόδου `output.file`, όπου θα γραφτεί το αποτέλεσμα της λειτουργίας `merge join`. Στη συνέχεια ανοίγω τα αρχεία εισόδου για ανάγνωση και το αρχείο εξόδου για εγγραφή. Τα `R_file`, `S_file` και `output` είναι αντικείμενα αρχείων. Δημιουργώ τα `R_reader` και `S_reader` που είναι αντικείμενα αναγνώσης csv για την αναάγνωση δεδομένων από τα αρχεία εισόδου και ένα αντικείμενο εγγραφής csv `writer` για την εγγραφή δεδομένων στο αρχείο εξόδου. Η `csv.reader()` χρησιμοποιείται για την ανάγνωση δεδομένων CSV από τα αρχεία εισόδου και η `csv.writer()` χρησιμοποιείται για την εγγραφή δεδομένων CSV στο αρχείο εξόδου. Έπειτα χρησιμοποιώ τη συνάρτηση `next()` για να διαβάσω τη πρώτη γραμμή δεδομένων από κάθε αρχείο εισόδου. Για όσο η `R_row` και `S_row` υπάρχουν(δηλαδή για όσο υπάρχουν δεδομένα προς ανάγνωση από ένα από τα αρχεία εισόδου), εάν οι τιμές των κλειδιών από τα `R_file` και `S_file` είναι ίσες γράφει μια νέα γραμμή στο αρχείο εξόδου με τη συνένωση των τρέχων γραμμών των `R_file` και `S_file` παραλείποντας την δεύτερη στήλη του `S_file`. Αφού γράψει τη γραμμή , προχωρά στην επόμενη γραμμή του `S_file` χρησιμοποιώντας `next(S_reader, None)`. Εάν η τιμή του κλειδιού από το `R_file` είναι μικρότερη από την τιμή του κλειδιού από το `S_file` τότε προχωράει στην επόμενη γραμμή του αρχείου `R_file` χρησιμοποιώντας `next(R_reader, None)`. Αλλιώς αν η τιμή του κλειδιού από το `R_file` είναι μεγαλύτερη από την τιμή του κλειδιού από το `S_file` τότε προχωράει στην επόμενη γραμμή του αρχείου `S_file` χρησιμοποιώντας `next(S_reader, None)`. Παρακάτω αφού βγεί από το `while` , εάν έχουν μείνει γραμμές στο `R_file` τις γράφει στο αρχείο εξόδου με κενά πεδία για τις στήλες από το `S_file`. Αντίστοιχα η επόμενη `while` εάν έχουν μείνει γραμμές στο `S_file` τις γράφει στο αρχείο εξόδου με κενά πεδία για τις στήλες από το `R_file`. Ουσιαστικά εκτελώ ένα `merge join` σε δυο αρχεία csv με βάση τα ακέραια κλειδιά που βρίσκονται στην πρώτη στήλη κάθε αρχείου και γράφω τα αποτελέσματα σε ένα άλλο αρχείο εξόδου.

Μέρος 3(composite query)

Αρχικά ορίζω τη συνάρτηση `my_composite_query` που λαμβάνει τρεις παραμέτρους: `R_file`, `S_file` και `output_file`. Αυτές οι παράμετροι αντιπροσωπεύουν τις διαδρομές προς τα δύο αρχεία CSV εισόδου `R_file` και `S_file` και το αρχείο CSV εξόδου `output.file`, όπου θα γραφτεί το αποτέλεσμα. Αρχικοποιώ ένα κενό λεξικό `sum_e_dictionary` για να αποθηκεύσω τα αθροίσματα των τιμών `E` για κάθε ξεχωριστή τιμή `A` που βρίσκεται στο αρχείο `S.csv`. Στη συνέχεια ανοίγω τα αρχεία εισόδου για ανάγνωση και το αρχείο εξόδου για εγγραφή. Αρχικοποιώ ένα csv αντικείμενο `writer` για εγγραφή στο αρχείο εξόδου. Οι γραμμές `R_line` και `line_s` αρχικοποιούνται διαβάζοντας την πρώτη γραμμή από το `R_file` και το `S_file` αντιστοίχα χρησιμοποιώντας την `readline()`. Αρχικοποιώ και μια κενή λίστα `toadd` για να αποθηκεύσω την τιμές από το `R_file` βάση της συνθήκης. Για όσο η `R_line` και `S_line` υπάρχουν(δηλαδή για όσο υπάρχουν δεδομένα προς ανάγνωση από ένα από τα αρχεία εισόδου) θα είναι μέσα στο βρόχο. Με τα `R_items = R_line.split(',')` και `S_items =`

`S_line.split(',')` χωρίζω τις γραμμες που διαβάζονται απο τα αρχεία εισόδου σε λίστες στοιχείων χρησιμοποιώντας την `split()`. Έπειτα εξάγω συγκεκριμένες τιμές απο τις `split lines`. Η μεταβλητή `a_r` αντιπροσωπευει την τιμη του πρώτου στοιχείου απο το `R_file` ως `int`. Η μεταβλητή `c` αντιπροσωπευει την τιμη του τρίτου στοιχείου απο το `R_file` ως `int`. Η μεταβλητή `a_s` αντιπροσωπευει την τιμη του δεύτερου στοιχείου απο το `S_file` ως `int`. Η μεταβλητή `e` αντιπροσωπευει την τιμη του τρίτου στοιχείου απο το `S_file` ως `int`. Παρακάτω ενημερώνω το `sum_e_dictionary`. Ελέγχω αν το κλειδί `a_s` υπάρχει ήδη στο λεξικό και αν όχι αρχικοποιώ την τιμή σε 0. Στη συνέχεια προσθέτω την τιμή του `e` στην υπάρχουσα τιμή που είναι associated με το κλειδι `a_s` στο λεξικό. Εάν η τιμή του `c` απο το `R_file` είναι ίση με 7 προσθέτω την τιμή της `a_r` στη λιστα `toadd`. Διαβάζω τις επόμενες γραμμές απο τα αρχεία εισόδου προχωρώντας έτσι στην επόμενη επανάληψη του `while`. Με τη `for` για κάθε στοιχείο της `toadd` ελέγχω αν το `x` υπάρχει ως κλειδί στο λεξικό. Εάν το `x` υπάρχει στο λεξικό γραφει μια νέα γραμμη στο αρχείο εξόδου που περιέχει την τιμή του `x` και το αθροισμα που σχετίζεται με αυτό στο λεξικό. Ουσιαστικά δηλαδή η συνάρτηση διαβάζει τιμές και απο τα δυο αρχεια, υπολογίζει αθροισματα με βάση τις συνθήκες και γράφει τα αποτελέσματα σε ένα αρχείο εξόδου.