

Διαχείριση Σύνθετων Δεδομένων

Assignment1

Φωτόπουλος Στέφανος ,4829

Ιωάννινα, 19/03/2023

Contents

Assignment1	1
1.Περιγραφή της εργασίας	3
Μέρος 1	3
Μέρος 2	5
Μέρος 3	6

1.Περιγραφή της εργασίας

Μέρος 1

Group-By With Aggregation:

```
# Stefanos Fotopoulos 4829
2 import sys
3 import csv
4 from collections import defaultdict
5
6
7 def merge_sort(array, attribute_index):
8     if len(array) <= 1:
9         return array
10
11     mid = len(array) // 2 # xwrizw ton pinika se 2 mikrotous
12     first_half = merge_sort(array[:mid], attribute_index) #anadromiki klisi tis sinartisis mexri na exei mono 1 stoixeio kratontas mono to attribute pou mas endiaferei
13     second_half = merge_sort(array[mid:], attribute_index) #anadromiki klisi tis sinartisis mexri na exei mono 1 stoixeio kratontas mono to attribute pou mas endiaferei
14     sorted_array = []
15     first_half_index = 0
16     second_half_index = 0
17
18     while first_half_index < len(first_half) and second_half_index < len(second_half):
19         if int(first_half[first_half_index][attribute_index]) < int(second_half[second_half_index][attribute_index]): #elegxw poia timi einai megalyteri apo tis 2 gia na tin eisagw ston n
20             sorted_array.append(first_half[first_half_index])
21             first_half_index += 1
22         else:
23             sorted_array.append(second_half[second_half_index])
24             second_half_index += 1
25
26     sorted_array.extend(first_half[first_half_index:]) #edw eisagoume ta enapomeinanta stoixeia ta opoia den mpikan ston sorted array
27     sorted_array.extend(second_half[second_half_index:]) #edw eisagoume ta enapomeinanta stoixeia ta opoia den mpikan ston sorted array
28     return sorted_array
```

Στη `merge_sort` ελέγχω κάθε φορά στην αρχή κατά πόσο το μήκος του πίνακά μου έχει φτάσει τη μονάδα καθώς αφού τον χωρίσω στη μέση σε δυο υποπίνακες θα την καλέσουν αναδρομικά. Με αυτή την αναδρομή επιτυγχάνω να σπάσω σε μικρότερα κομμάτια τον αρχικό μου πίνακα, έτσι ώστε να γίνει ευκολότερη η ταξινόμισή του. Αφού τελειώσω με τις αναδρομές θα χρειαστεί να ανασυνθέσω τον παλιό μου πίνακα για αυτό τον λόγο υπάρχει η `while`. Καθώς σε κάθε επανάληψη κοιτάει ποια από τις δυο τιμές των νέων πινάκων είναι μεγαλύτερη για να την εισάγει στον `sorted` μεγαλώνοντας το `index` αυτού του πίνακα κατά ένα για να ελέγξω την επόμενη τιμή του. Τέλος, υπάρχει περίπτωση να μείνουν κάποια στοιχεία που δεν έχουν εισαχθεί στον `sorted_array` για αυτό χρησιμοποιώ τα `extend`.

```

31 def group_by_aggregation(sorted_data, group_by_attr, aggregation_attr, function):
32     data = defaultdict(list)
33
34     for row in sorted_data: #diatrexei kathe grammi tou sorted_data pou phre apo tin merge sort
35         group_by = row[group_by_attr]
36         aggregation_value = row[aggregation_attr]
37         data[group_by].append(aggregation_value) #prostheti tin aggregation_value sto leksiko analoga me to ekastote group by
38
39     with open('01.csv', 'w', newline='') as output:
40         writer = csv.writer(output)
41         for group_by, aggregation_values in data.items(): #anatrexei gia to sygkekrimeno kleidi(group_by) stis times(aggregation_values) pou exei to leksiko
42             if function == 'sum':
43                 result = sum(map(int, aggregation_values)) #xrelazetai map gia na paroume ton akeraio
44             elif function == 'min':
45                 result = min(map(int, aggregation_values)) #xrelazetai map gia na paroume ton akeraio
46             elif function == 'max':
47                 result = max(map(int, aggregation_values)) #xrelazetai map gia na paroume ton akeraio
48             writer.writerow([group_by, result])
49

```

Για την group_by_aggregation δημιουργώ ένα λεξικό που θα με βοηθήσει να βρω τις τιμές για κάθε group_by_attribute. Αρχικά, διατρέχω τα δεδομένα που έγιναν sort από την προηγούμενη συνάρτηση και αρχικοποιώ τις μεταβλητές μου. Προσθέτω στο λεξικό μου την τιμή για το εκάστοτε κλειδί σύμφωνα με το group_by. Δημιουργώ ένα νέο αρχείο για εγγραφή όπως και έναν writer για αυτό το αρχείο. Με μια for loop ανατρέχω στο λεξικό μου για κάθε κλειδί της group_by και ανάλογα με το επιθυμητό Function δίνω το αποτέλεσμα το οποίο γράφω στο νέο αρχείο δίπλα από το group_by_attribute.

```

51 if __name__ == "__main__":
52     if len(sys.argv) != 5:
53         print("Warning follow the pattern: python program.py <filename> <group_by_attribute> <aggregation_attribute> <function>")
54         sys.exit(1)
55
56     filename = sys.argv[1]
57     group_by_attribute = int(sys.argv[2])
58     aggregation_attribute = int(sys.argv[3])
59     aggregation_function = sys.argv[4]
60
61     with open(filename, newline='') as file:
62         reader = csv.reader(file)
63         data = list(reader)
64
65     sorted_data = merge_sort(data, group_by_attribute)
66
67     group_by_aggregation(sorted_data, group_by_attribute, aggregation_attribute, aggregation_function)
68

```

Στην main δείχνω στον χρήστη ποια είναι η σωστή σειρά των arguments που πρέπει να γράψει στην κονσόλα για να τρέξει το πρόγραμμα. Έπειτα, ανοίγω το αρχείο μου δημιουργώ έναν αναγνώστη για το αρχείο αυτό και διαβάζω τα δεδομένα από το αρχείο ως λίστες. Έτσι πλέον έχω μια λίστα από λίστες, όπου κάθε λίστα περιέχει τα δεδομένα από μια γραμμή του αρχικού csv αρχείου.

!!!!Προσοχή: Παρατήρησα ότι για να αλλάξουν τα δεδομένα στο νέο αρχείο θα χρειαστεί πρώτα να κλείσετε το παλιό(O1.csv) δηλαδή αν τρέξετε το αρχείο για max το ανοίξετε στο Excel και ενώ το έχετε ανοιχτό τρέξετε το πρόγραμμα για min θα σας πετάξει PermissionError: [Errno 13] Permission denied: '01.csv'. Εάν κλείσετε το αρχείο(O1.csv για το Max) θα αλλάξει το 01.csv και θα δείτε το 01.csv με τις αλλαγμένες πλέον τιμές για min.

Μέρος 2

Merge join:

```
1 #Stefanos Fotopoulos 4829
2 import csv
3
4 def merge_join(R_csv, S_csv, output):
5     with open(R_csv, newline='') as R, open(S_csv, newline='') as S, open(output, 'w', newline='') as output_csv:
6         #Dimiourgia readers/writer gia na mporesw na diatreksw to arxeio kai na grapsw ta apotelesmata sto output
7         reader_R = csv.reader(R)
8         reader_S = csv.reader(S)
9         writer = csv.writer(output_csv)
10
11         for row_R in reader_R:
12             S.seek(0) #Epanaferw ton file pointer stin arxi tou S
13             for row_S in reader_S:
14                 if row_R[0] == row_S[1]: #R(ABC) S(DAE) elegxw tin prwth stili tou R me tin mesaia tou S
15                     writer.writerow(row_R + row_S[:1] + row_S[2:]) #Krataw mono ABC apo R kai DE apo S
16
17 if __name__ == "__main__":
18     merge_join('R.csv', 'S.csv', '02.csv')
19
```

Με την βοήθεια της with ανοίγω τα αρχεία .csv για ανάγνωση και εγγραφή και δημιουργώ readers και writer έτσι ώστε να διατρέξω το αρχείο και να γράψω μέσω του writer τα αποτελέσματα στο αρχείο output. Με το for loop που έπεται διατρέχω την κάθε γραμμή της πρώτης στήλης του R με τις εγγραφές της μεσαίας στήλης του S και εάν αυτές οι τιμές είναι ίσες κρατάω το αποτέλεσμα(ABCDE) στο αρχείο output. Τέλος, θα χρειαστεί στην αρχή κάθε 'τρέξιματος' να επαναφέρουμε τον file pointer του S στο 0 για να συγκριθεί με την επόμενη τιμή του R.

Μέρος 3

Composite Query:

```
1 #Stefanos Fotopoulos 4829
2 import csv
3
4 def composite_query(R_csv, S_csv, output):
5
6     with open(R_csv, newline='') as R, open(S_csv, newline='') as S, open(output, 'w', newline='') as output_csv:
7         # Dimiourgia readers/writer gia na mporesw na diatreksw to arxeio kai na grapsw ta apotelesmata sto output
8         reader_S = csv.reader(S)
9         reader_R = csv.reader(R)
10        writer = csv.writer(output_csv)
11
12        for row_R in reader_R:
13            R_A = row_R[0]
14            C = int(row_R[2])
15            sum = 0
16
17            if C == 7:
18                S.seek(0) # Epanaferw ton file pointer stin arxi tou S
19                for row_S in reader_S:
20                    S_A = row_S[1]
21                    E = int(row_S[2])
22                    if R_A == S_A:
23                        sum += E
24                if sum > 0: # Xwris auth thn synthiki krataei kai tis times pou exoun mono R.C = 7
25                    writer.writerow([R_A, sum])
26
27
28 if __name__ == "__main__":
29     composite_query('R.csv', 'S.csv', '03.csv')
```

Με την βοήθεια της with ανοίγω τα αρχεία .csv για ανάγνωση και εγγραφή και δημιουργώ readers και writer έτσι ώστε να διατρέξω το αρχείο και να γράψω μέσω του writer τα αποτελέσματα στο αρχείο output. Αρχικοποιώ τα R_A, C, sum (μεταβλητή στην οποία θα κρατήσω το άθροισμα) σύμφωνα με τα δεδομένα της εκφώνησης δηλαδή το R_A είναι η 1^η στήλη του R και το C η 3^η όπου συγκρατώ τον integer. Ελέγχω πρώτα για το αν η τιμή του C είναι 7 και εάν ισχύει δίνω τιμές στα S_A και E. Τέλος, όταν ισχύει η ισότητα προσθέτω στο sum την τιμή του E, θα χρειαστεί στην αρχή κάθε 'τρέξιματος' να επαναφέρουμε τον file pointer του S στο 0 για να συγκριθεί με την επόμενη τιμή του R.