

Διαχείριση Σύνθετων Δεδομένων

Assignment2

Φωτόπουλος Στέφανος ,4829

Ιωάννινα,03/05/2023

Contents

Assignment2.....	1
1.Περιγραφή της εργασίας	3
Μέρος 1	3
Μέρος 2	10

1.Περιγραφή της εργασίας

Μέρος 1

```
6 def read_data(filename):
7     with open(filename) as f:
8         num_points = int(f.readline())
9         print("Total Coordinates: ", num_points, "\n")
10        points = []
11        for i, line in enumerate(f.readlines(), start=1):
12            point = [i] + list(map(float, line.strip().split()))
13            points.append(point)
14        return points, num_points
15
16
17 def sort_points_by_x(points):
18     sorted_points = sorted(points, key=lambda p: p[1]) #sort by x
19     return sorted_points
20
21
22 def sort_points_by_y(points, num_points):
23     num_per_node = math.floor(1024 / 20) #51 στοιχεία θα χει το κάθε φύλλο
24     num_nodes = math.ceil(num_points / num_per_node) #Synolika 1020 fylla
25     square_root = math.ceil(math.sqrt(num_nodes))
26     vertical_2d_lines = square_root * num_per_node
27     sorted_points = []
28
29     #Xwrizw ton pinaka se 1632 tmimata kai ta taksinomw kata y
30     for i in range(0, len(points), vertical_2d_lines):
31         segment = points[i:i + vertical_2d_lines]
32         sorted_segment = sorted(segment, key=lambda p: p[2]) #sort by y
33         sorted_points.extend(sorted_segment)
34
35     return sorted_points
```

Αρχικά, στην read_data διαβάζω την πρώτη γραμμή και την αποθηκεύω στην μεταβλητή num_points τυπώνοντας το σύνολο των συντεταγμένων. Αποθηκεύω σε μια λίστα τις υπόλοιπες τιμές και τις επιστρέφω.

Sort_points_by_x: Ταξινομώ την λίστα μου με τις τιμές του x και την επιστρέφω

Sort_points_by_y: Βρίσκω πόσα στοιχεία θα χει κάθε φύλλο (1024B/20B) (Παίρνω το floor γιατί θα χρησιμοποιηθούν τα 1020B από τα συνολικά 1024 δεν μπορώ να υπερβώ αυτόν τον αριθμό), τα συνολικά φύλλα μου είναι 1020 και μετά χρησιμοποιώ τον τύπο του μαθήματος $\sqrt{N} * M$ όπου N ο συνολικός αριθμός των φύλλων και M των συνολικών στοιχείων του εκάστοτε φύλλου. Ταξινομώ τις συντεταγμένες μου πρώτα κατά x και έπειτα ανά 1632 κατά y και επιστρέφω την τελική μου λίστα.

```

38 def minimum_bounding_rectangle(data, flag):
39     mbr = []
40     if flag == 0:
41         flattened_data = data
42     else:
43         flattened_data = [item for sublist in flattened_data for item in sublist]
44     else:
45         flattened_data = data
46     for i in range(0, len(flattened_data)):
47         min_X = min(item[1] for item in flattened_data[i])
48         max_X = max(item[1] for item in flattened_data[i])
49         min_Y = min(item[2] for item in flattened_data[i])
50         max_Y = max(item[2] for item in flattened_data[i])
51         mbr.append([min_X, min_Y, max_X, max_Y])
52     return mbr

```

Στην συνάρτηση `minimum_bounding_rectangle` : Έχω δύο ορίσματα τα δεδομένα μου και το `flag`(μου δείχνει αν το καλεί κόμβος ή φύλλο). Κύρια εργασία της είναι να βλέπει κάθε γραμμή των δεδομένων και να επιστρέφει σε μια λίστα τις `min/max` τιμές των αξόνων `x,y`.

```

55 def find_best_mbr(data, total_areas):
56     areas_list = []
57     flattened_data = data
58     counter = 0
59     flattened_data = [item for sublist in flattened_data for item in sublist]
60     index_list = []
61     mbr_list = []
62     for i in range(0, len(flattened_data)):
63         counter = 0
64         for item in flattened_data[i][0]:
65             min_X = item[0]
66             min_Y = item[1]
67             max_X = item[2]
68             max_Y = item[3]
69             width = max_X - min_X
70             height = max_Y - min_Y
71             temp_area = width * height
72             total_areas.append(temp_area)
73             if counter == 0:
74                 areas_list.append(temp_area)
75                 index_list.append(counter)
76             if temp_area > areas_list[i]:
77                 areas_list.pop(i)
78                 index_list.pop(i)
79                 areas_list.append(temp_area)
80                 index_list.append(counter)
81             counter += 1
82
83     for i in range(0, len(index_list)):
84         index = index_list[i]
85         mbr_list.append(flattened_data[i][0][index])
86     return mbr_list
87

```

Η `find_best_mbr` βρίσκει το καλύτερο `mbr` μέσω του εμβαδού για κάθε κόμβο και τα επιστρέφει σε μια λίστα. Κρατώ σε μια `index_list` τα `Indexes` με τα καλύτερα εμβαδά έτσι ώστε στο τελευταίο `for loop` να ανατρέξω στην αρχική μου λίστα και να πάρω μόνο

τα mbrs που με ενδιαφέρουν από κάθε κόμβο. Χρησιμοποιώ τον counter για να με βοηθήσει να γνωρίζω σε ποια θέση του αρχικού μου πίνακα βρίσκεται το καλύτερο Mbr, στο πρώτο run θεωρώ ως καλύτερο εμβαδό αυτό της πρώτης θέσης και στην συνέχεια βλέπω αν τα επόμενα είναι μεγαλύτερα από αυτό που ήδη έχω σε περίπτωση που ισχύει κάνω pop την προηγούμενη τιμή και εισάγω την νέα. Τέλος, κρατάω σε μια λίστα όλες τις τιμές των εμβαδών των κόμβων(total_areas).

```
89 def average_mbr_area(nodes, num_nodes, num_of_last_nodes):
90     total_area = 0
91     for i in range(0, len(nodes)):
92         total_area += nodes[i]
93     total_nodes = (num_nodes-1) * 28 + num_of_last_nodes
94     result = total_area / total_nodes
95     return result
96
```

Average_mbr_area: Στην λίστα nodes έχω όλες τις τιμές των εμβαδών από το επίπεδο στο οποίο βρίσκομαι, num_nodes είναι ο αριθμός των συνολικών κόμβων σε αυτό το επίπεδο, num_of_last_nodes είναι ο αριθμός των στοιχείων του τελευταίου κόμβου. Για την εύρεση του average_mbr προσθέτω όλα τα εμβαδά και διαιρώ τους συνολικούς κόμβους μείον έναν(που είναι ο τελευταίος) επί 28 αφού έχει 28 στοιχεία ο κάθε ολοκληρωμένος κόμβος συν τα στοιχεία του τελευταίου.

```

98 def build_leaf(coordinates, num_points, level):
99     num_per_node = math.floor(1024 / 20)
100     num_nodes = math.ceil(num_points / num_per_node)
101     leafs = []
102     id = 0
103     def sort_tile_recursive(coords, node_id):
104         nonlocal id
105         if not coords:
106             return []
107         leafs_for_node = []
108         while coords:
109             tile = coords[:num_per_node]
110             sorted_coords_y = sorted(tile, key=lambda x: x[2])
111
112             leaf_bytes = 0
113             temp_leaf = []
114             for coord in sorted_coords_y:
115                 if leaf_bytes + 20 > 1024:
116                     leafs_for_node.append(temp_leaf)
117                     leaf_bytes = 0
118                     temp_leaf = []
119                 temp_leaf.append(coord)
120                 leaf_bytes += 20
121
122             if temp_leaf:
123                 leafs_for_node.append(temp_leaf)
124                 temp_leaf = []
125             coords = coords[num_per_node:]
126             id = node_id + 1
127             return leafs_for_node
128     leafs = sort_tile_recursive(coordinates, id)
129     level += 1
130     return leafs

```

build_leaf: λαμβάνει ως είσοδο μια λίστα `coordinates` που περιέχει συντεταγμένες σημείων, τον αριθμό των σημείων `num_points` και το επίπεδο του κόμβου `level`. Υπολογίζει τον αριθμό των κόμβων που θα δημιουργηθούν, με βάση τον μέγιστο αριθμό σημείων που μπορεί να περιέχει ένας κόμβος. Η εσωτερική συνάρτηση `sort_tile_recursive`, η οποία ταξινομεί τις συντεταγμένες των σημείων. Ο στόχος είναι να οργανωθούν τα σημεία σε τετράγωνα (tiles) με βάση τον άξονα `y`. Αρχικά, δημιουργούνται tiles μεγέθους `num_per_node` (στην περίπτωσή μας 20) και ταξινομούνται βάσει του `y`. Έπειτα, τα tiles που δημιουργούνται συσσωρεύονται σε μια λίστα `leafs_for_node`. Η διαδικασία αυτή επαναλαμβάνεται για όλα τα tiles που περιέχονται στη λίστα `coordinates`. Τέλος, καλείται η `sort_tile_recursive` με τις συντεταγμένες `coordinates` και επιστρέφει τη λίστα των tiles `leafs_for_node`, η οποία αποθηκεύεται στη μεταβλητή `leafs`. Το επίπεδο `level` αυξάνεται κατά 1 και επιστρέφεται η λίστα των `leafs`.

```

133 def build_rtree(leafs_list, R_tree, total_areas, temp_mbr, average_list, flag, num_points):
134     total_bytes = 1024
135     num_per_node = math.floor(total_bytes / 20) # ta shmeia poy xwraei kathe fyllo (51)
136     num_nodes = math.ceil(num_points / num_per_node) # synolika fyllo (1020)
137     leafs_per_node = math.floor(total_bytes / 36) # fyllo ana komvo (28)
138     total_nodes = math.ceil(num_nodes / leafs_per_node) # synolo komvwn (37)
139     mbr_bytes = 32 + 4
140     original_list = leafs_list
141     first_len = len(leafs_list)
142     average_list.append(first_len)
143     leafs = leafs_list
144     mbr = []
145     counter_leaf = 0
146     counter_node = 0
147     node_id = 0
148     last_mbr = 1
149
150
151     if len(leafs) == num_nodes:
152         flag = 0
153         for i in range(1, len(leafs)+1):
154             if i % leafs_per_node == 0 or i == num_nodes:
155                 mbr = []
156
157                 if counter_node + mbr_bytes > total_bytes:
158                     node_coords = []
159                     node_coords.append(original_list[counter_node:first_len])
160                     mbr.append(minimum_bounding_rectangle(node_coords, flag))
161                     node_id += 1
162                     leafs.append(mbr)
163                     R_tree.append(mbr)
164                     temp_mbr.append(len(node_coords[0]))
165                     counter_node += leafs_per_node
166                 else:
167                     node_coords = []
168                     node_coords.append(leafs[counter_node:counter_node + leafs_per_node])
169                     mbr.append(minimum_bounding_rectangle(node_coords, flag))
170                     last_mbr += 1
171                     node_id += 1
172                     leafs.append(mbr)
173                     R_tree.append(mbr)
174
175                     counter_node += leafs_per_node
176                 counter_leaf += num_per_node
177             else:
178                 flag = 1
179                 for i in range(1, len(leafs_list)+1):
180                     if i % leafs_per_node == 0 or i == first_len:
181                         mbr = []
182                         if i == first_len:
183                             node_coords = []
184                             node_coords.append(original_list[counter_node:first_len])
185                             mbr_list = find_best_mbr(node_coords, total_areas)
186                             mbr.append(mbr_list)
187                             last_mbr += 1
188                             node_id += 1
189                             leafs.append(mbr)
190                             if first_len != 1:
191                                 R_tree.append(mbr)
192                                 temp_mbr.append(len(node_coords[0]))
193                                 counter_node += leafs_per_node
194                         else:
195                             node_coords = []
196                             node_coords.append(original_list[counter_node:counter_node + leafs_per_node])
197                             mbr_list = find_best_mbr(node_coords, total_areas)
198                             mbr.append(mbr_list)
199                             last_mbr += 1
200                             node_id += 1
201                             leafs.append(mbr)
202                             R_tree.append(mbr)
203                             counter_node += leafs_per_node
204                         counter_leaf += num_per_node
205                 if flag == 1:

```

build_rtree: Όπως και στις προηγούμενες συναρτήσεις αρχικοποιώ τις μεταβλητές που χρειάζομαι μπορούμε να πούμε ότι χωρίζεται σε δυο μέρη το ένα είναι όταν παίρνει όρισμα τα φύλλα του δέντρου και μετά είναι οι υπόλοιποι κόμβοι. Για την πρώτη περίπτωση διατρέχω όλη την λίστα και ανά 28 στοιχεία βρίσκω το mbr κάθε γραμμής και τα αποτελέσματα προστίθενται στο R_tree. Όταν βρίσκομαι στην δεύτερη περίπτωση θέτω το flag 1 βρίσκω το καλύτερο mbr κάθε κόμβου και το εισάγω στο R_tree. Για τις τελευταίους κόμβους τόσο στα φύλλα όσο και στους εσωτερικούς κόμβους η λογική είναι η ίδια απλά κρατάω και τον αριθμό των συνολικών στοιχείων των φύλλων σε μια λίστα temp_mbr ομοίως και για τους εσωτερικούς κόμβους.

```

205         if flag == 1:
206             num_of_last_nodes = temp_mbr[0]
207             avg_mbr = average_mbr_area(total_areas, first_len, num_of_last_nodes)
208             temp_mbr.pop(0)
209             average_list.append(avg_mbr)
210             total_areas = []
211             if first_len == 1:
212                 return R_tree
213
214             last_level_list = leaves[first_len:]
215
216             if first_len != 1:
217                 build_rtree(last_level_list, R_tree, total_areas, temp_mbr, average_list, flag, num_points)
218
219

```

Όταν η τιμή flag γίνει 1 σημαίνει ότι έχω προχωρήσει από τα φύλλα και βρίσκομαι σε κάποιον εσωτερικό κόμβο παίρνω την πρώτη τιμή της λίστας temp_mbr η οποία είναι το πλήθος στοιχείων του τελευταίου κόμβου αυτού του επιπέδου, αυτός ο αριθμός θα με βοηθήσει στον υπολογισμό του average_mbr. Ελέγχω εάν βρίσκομαι στην ρίζα του δέντρου και αν όχι ξανακαλώ την build_rtree με καινούρια λίστα αρχίζοντας πλέον από το σημείο που βρίσκεται ο τελευταίος κόμβος του προηγούμενου επιπέδου.

```

221 def print_tree_statistics(R_tree, avg_list, free_stats, level):
222     Stats = ""
223     Tree_Stats += "Root id is: " + str(len(R_tree)-1) + "\n"
224     Tree_Stats += "          Leaf Node          \n"
225     indices_to_move = [0, 1, 3, 5]
226     num_of_nodes = [avg_list.pop(i) for i in reversed(indices_to_move)]
227     for i in range(0, num_of_nodes[-1]):
228         Stats = str(i) + ", " + str(len(R_tree[i])) + ", " + str(0)
229         coords = ""
230         for j in range(0, len(R_tree[i])):
231             coords += ", " + "(" + str(R_tree[i][j][0]) + ", " + "(" + str(R_tree[i][j][1]) + ", " + str(R_tree[i][j][2]) + ")" + ")"
232         Tree_Stats += Stats + coords + "\n"
233
234     print("Total number of leaves: ", num_of_nodes[-1], " found in level", level)
235     print("Average MBR for level ", level, "is: 0.0")
236     first_len_of_avg_list = len(avg_list)
237     last_num = num_of_nodes[-1]
238     num_of_nodes.pop(-1)
239     level += 1
240     record_id = 0
241     while len(num_of_nodes) != 0:
242         Tree_Stats += "          Inner Node at Level " + str(level) + "          \n"
243         for i in range(last_num, last_num + num_of_nodes[-1]):
244
245             Stats = str(i) + ", " + str(len(R_tree[i][0])) + ", " + str(1)
246             coords = ""
247             for j in range(0, len(R_tree[i][0])):
248                 coords += ", " + "(" + str(record_id) + ", " + "(" + str(R_tree[i][0][j][0]) + ", " + str(R_tree[i][0][j][1]) + ", " + str(R_tree[i][0][j][2]) + ")" + ")"
249                 record_id += 1
250             Tree_Stats += Stats + coords + "\n"
251     print("Total number of internal nodes: ", num_of_nodes[-1], " found in level", level)
252     print("Average MBR for level ", level, "is: ", avg_list[0])
253     avg_list.pop(0)
254     level += 1
255     last_num += num_of_nodes[-1]
256     num_of_nodes.pop(-1)
257     print("Height of R_Tree is: ", first_len_of_avg_list + 1)
258
259     return Tree_Stats

```

Print_tree_statistics: Παίρνει ως ορίσματα το R_tree που είναι η λίστα με όλες τις τιμές που θέλουμε, το avg_list έχει μέσα το πλήθος κόμβων κάθε επιπέδου όπως και τα average_mbr του κάθε επιπέδου. Στην πρώτη γραμμή του αρχείου τυπώνω το Id της ρίζας και μετά το Leaf node με όλα τα φύλλα, στη συνέχεια όλους τους εσωτερικούς κόμβους ανά επίπεδο μαζί με το average mbr του κάθε επιπέδου. Τέλος, τυπώνω το ύψος του δέντρου που είναι οι τιμές του average mbr που έχω για κάθε επίπεδο προσθέτοντας ένα γιατί στο επίπεδο των φύλλων το average mbr είναι 0 οπότε δεν έχω κάποια τιμή.


```

262 def main(output_filename):
263     input_filename = "Beijing_restuarants.txt"
264     R_tree = []
265     total_areas = []
266     temp_mbr = []
267     average_list = []
268     Tree_Stats = ""
269     level = 0
270     flag = 0
271     coordinates, num_points = read_data(input_filename)
272     sorted_coords_x = sort_points_by_x(coordinates)
273     sorted_coords = sort_points_by_y(sorted_coords_x, num_points)
274     root = build_leaf(sorted_coords, num_points, level)
275
276     R_tree += root
277     nodes = build_rtree(root, R_tree, total_areas, temp_mbr, average_list, flag, num_points)
278     Tree_Stats = print_tree_statistics(R_tree, average_list, Tree_Stats, level)
279
280
281     with open(output_filename, "w") as file:
282         file.write(Tree_Stats)
283
284
285 if __name__ == "__main__":
286     if len(sys.argv) != 2:
287         print("Usage: python script_name.py output_filename")
288         sys.exit(1)
289
290     output_filename = sys.argv[1]
291     main(output_filename)
292

```

Έχω δημιουργήσει μια συνάρτηση Main με όρισμα το output αρχείο ουσιαστικά μέσα σε αυτήν αρχικοποίησα όλες τις global μεταβλητές που χρησιμοποιώ στις παραπάνω συναρτήσεις. Η διαδικασία της Main είναι να διαβάσει το αρχείο, να κάνει sort κατά x και μετά κατά y, να δημιουργήσει τα φύλλα και έπειτα τους εσωτερικούς κόμβους που τα δεδομένα αυτά θα εισαχθούν στο R_tree, με το R_tree θα παράξει τα στατιστικά του δέντρου, τα οποία θα τα εισάγει στο Output αρχείο.

Τέλος, η κανονική main του προγράμματος ελέγχει κατά πόσο ο χρήστης έχει καλέσει σωστά το πρόγραμμα στο τερματικό και καλεί την συνάρτηση main.

Μέρος 2

```
1 # Stefanos Fotopoulos 4829
2 import math
3
4 def read_tree_from_file(file_path):
5     result_list = []
6
7     with open(file_path, 'r') as file:
8         for line in file:
9             if not (line.startswith("Root id") or "----- Inner Node at Level 3 -----" in line or "----- Leaf Node -----" in line or "----- Inner Node at Level 1 -----" in line):
10                 result_list.append(line.strip())
11
12     return result_list
13
14
15 def mindist(point, mbr):
16     total = 0
17     for i in range(len(point)):
18         if point[i] < mbr[i]:
19             total += (mbr[i] - point[i]) ** 2
20         elif point[i] > mbr[i + len(point)]:
21             total += (point[i] - mbr[i + len(point)]) ** 2
22         else:
23             total += (point[i] - mbr[i]) ** 2
24     return math.sqrt(total)
25
26
27 if __name__ == "__main__":
28     import sys
29     if len(sys.argv) != 4:
30         print("Usage: python script_name.py rtree_file x,y k")
31         sys.exit(1)
32     tree_file = sys.argv[1]
33     q_coordinates = tuple(map(float, sys.argv[2].split(',')))
34     k = int(sys.argv[3])
35
36     tree = read_tree_from_file(tree_file)
```

Στο μέρος 2 έχω δυο συναρτήσεις την `read_tree_from_file` στην οποία εισάγω το αρχείο που έχει το δέντρο και μου επιστρέφει μια λίστα με tuples χωρίς τις περιττές γραμμές που έχω βάλει για να ξεχωρίζω κάθε επίπεδο.

Την `Mindist` η οποία παίρνει ένα σημείο(x,y) και το `mbr` ελέγχει κατά πόσο η τιμή του σημείου είναι μικρότερη από το Minimum `mbr` ή μεγαλύτερη από το maximum `mbr` και επιστρέφει την ρίζα της πρόσθεσης αυτών των δυο τετραγώνων(x και y).

Δυστυχώς δεν μπόρεσα να καταφέρω να εισάγω τα `mbr` από την λίστα μου οπότε το πρόγραμμα μου αυτή την στιγμή επιστρέφει σε μια τιμή `tree` αυτή την λίστα με τα tuples.