

Πολυδιάστατες Δομές Δεδομένων και Υπολογιστική Γεωμετρία
Εργαστηριακή Άσκηση 2022



<u>Ονοματεπώνυμο</u>	<u>ΑΜ</u>	<u>Έτος</u>	<u>e-Mail</u>
1. Χαριλάου Εφραίμ	1056638	6 ^ο	st1056638@ceid.upatras.gr
2. Δημητρίου Χρίστος	1056643	6 ^ο	st1056643@ceid.upatras.gr
3. Χατζή Βασιλική	1064359	6 ^ο	st1064359@ceid.upatras.gr
4. Μαρινόπουλος Γιώργος	1059573	6 ^ο	st1059573@ceid.upatras.gr
5. Κόλλιας Ιωάννης	1064886	6 ^ο	st1064886@ceid.upatras.gr

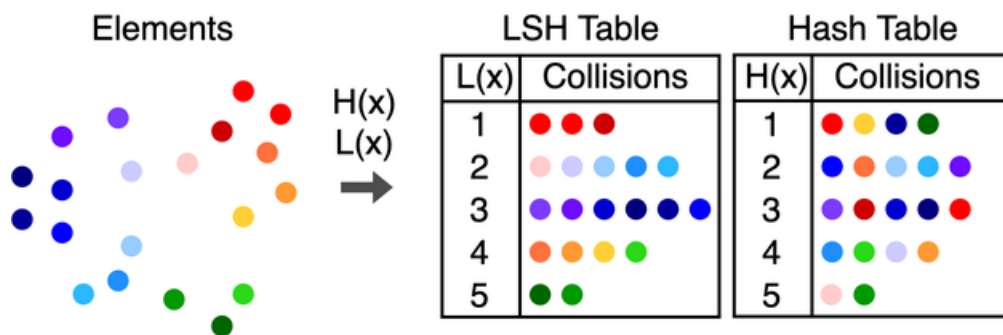
Περιεχόμενα:

LSH	3
Quad Tree.....	4
QTree.py	5
Local_QTree.py:.....	9
run_QTree.py	10
names.py:.....	10
Αποτελέσματα:	11
Range Search Tree.....	13
range_tree.py:.....	14
menu.py:	16
Αποτελέσματα:	17
R – Tree	20
Rtree.py:	21
Αποτελέσματα:	22
K-D Tree:	25
Kd-tree.py:	26
Αποτελέσματα:	27

LSH

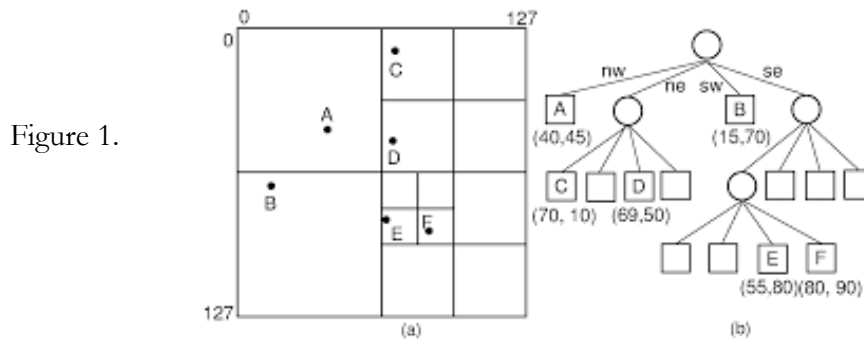
Η αναζήτηση πλησιέστερου γείτονα σε πολυδιάστατο χώρο είναι ένα σημαντικό πρόβλημα με πολλές εφαρμογές. Μια καλή λύση, όσον αφορά τη βάση δεδομένων, πρέπει να έχει δύο ιδιότητες: να μπορεί να ενσωματωθεί εύκολα σε μια σχεσιακή βάση δεδομένων και το κόστος αναζήτησης θα πρέπει να αυξάνεται υπογραμμικά σε σχέση με το μέγεθος του συνόλου δεδομένων, ανεξάρτητα από τις κατανομές δεδομένων και ερώτησης. Η μέθοδος κερματισμού με ευαισθησία τοπικότητας (Locality Sensitive Hashing - LSH) είναι μια γνωστή μέθοδος που ικανοποιεί και τις δύο προϋποθέσεις, όμως, οι τρέχουσες υλοποιήσεις της είτε έχουν υψηλό κόστος χώρου (απαιτούν δηλαδή αυξημένο αποθηκευτικό χώρο) και ερώτησης, είτε δεν διατηρούν τη θεωρητική της εγγύηση ποιότητας αποτελεσμάτων.

Ο αλγόριθμος LSH αποτελεί την βάση για περίπλοκα συστήματα καθώς έχει την δυνατότητα να συμπιέζει τις εγγραφές δημιουργώντας "signatures", ακολουθίες ακεραίων αριθμών έτσι ώστε να έχουμε τη δυνατότητα να συγκρίνουμε τίτλους χωρίς να χρειάζεται να διαχειριστούμε ολόκληρα σύνολα λέξεων. Εκτελεί αναζητήσεις πλησιέστερου γείτονα βασιζόμενος στην απλή ιδέα της "ομοιότητας". Λέμε ότι δύο αντικείμενα είναι όμοια όταν η τομή των συνόλων τους είναι επαρκώς μεγάλη.

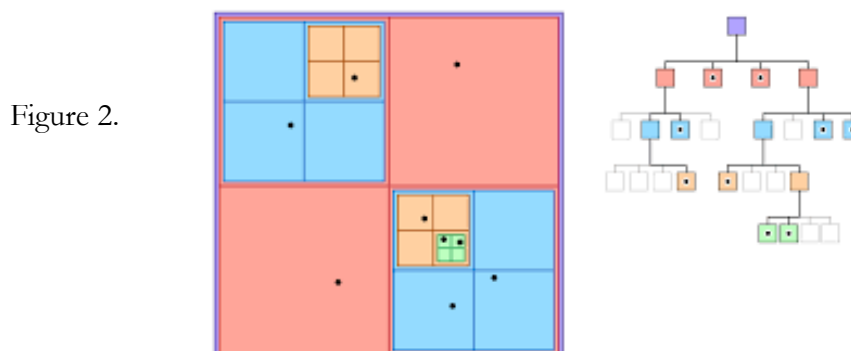


Quad Tree

Το Quadtree (Τετράδεντρο) είναι μια δεντρική δομή δεδομένων στην οποία κάθε εσωτερικός κόμβος έχει ακριβώς τέσσερα παιδιά. Συνήθως χρησιμοποιείται για την κατάρτιση ενός δισδιάστατου χώρου σε ορθογώνιες περιοχές. Κάθε κόμβος φύλλου στο Quadtree αντιπροσωπεύει μια περιοχή και η τιμή που αποθηκεύεται σε αυτόν τον κόμβο φύλλου είναι τα δεδομένα που σχετίζονται με την περιοχή αυτή. Το δέντρο μπορεί να χρησιμοποιηθεί για διάφορες λειτουργίες, όπως η αναζήτηση ενός συγκεκριμένου σημείου ή περιοχής, η εισαγωγή ενός νέου σημείου ή περιοχής και η διαγραφή ενός σημείου ή περιοχής. Η χρήση ενός Quadtree έχει το πλεονέκτημα της σημαντικής βελτίωσης στην αποτελεσματικότητα αυτών των λειτουργιών σε περιπτώσεις όπου τα δεδομένα είναι κατανεμημένα σε ένα μεγάλο δισδιάστατο χώρο.



Στην ενότητα αυτή της εργασίας, το κέντρο ενδιαφέροντος στρέφεται στα 2D Quad-Tree. Η διαδικασία της υλοποίησης και αναζήτησης έχει ως εξής: Ορίζεται ένα δισδιάστατο Ευκλείδειο επίπεδο, στο οποίο εξετάζεται ένα κανονικοποιημένο τετράγωνο. Στη συνέχεια, δημιουργείται μια ρίζα για το δέντρο, η οποία αντιπροσωπεύει το συνολικό τετράγωνο στο οποίο γίνεται η αναζήτηση. Κάθε κόμβος (συμπεριλαμβανομένης και της ρίζας) αντιπροσωπεύει μία περιοχή του τετραγώνου. Η διαδικασία που ακολουθεί ο αλγόριθμος είναι η επαναληπτική υποδιάρθρωση του κάθε κόμβου σε 4 τεταρτημόρια έως ότου βρεθούν όλα τα στοιχεία αναζήτησης και τοποθετηθούν στο δέντρο, σύμφωνα με τους κανόνες που το διέπουν.



Η εκπόνηση αυτού του μέρους της εργασίας έγινε στο περιβάλλον του Visual Studio Code Editor, καθώς και οι βιβλιοθήκες που χρησιμοποιήθηκαν παρατίθενται ακολούθως:

1. import math
2. from collections import deque
3. import time
4. import re
5. import pandas as pd

QTree.py

Στο αρχείο αυτό περιέχονται όλες οι απαραίτητες μέθοδοι, των οποίων η κλήση αφορά στη δημιουργία ενός τετραγώνου αναζήτησης, στην τοποθέτηση των στοιχείων στο κανονικοποιημένο επίπεδο, στη διαδικασία των υποδιαιρέσεων των τετραγώνων καθώς και στην εισαγωγή των κόμβων στο δέντρο.

Class Point(object):

Οι συντεταγμένες που θα χαρακτηρίζει κάθε σημείο στο χώρο με την επιλογή να του δοθούν χαρακτηριστικά δεδομένα.

Class BoundingBox():

Η αναπαράσταση γίνεται σε ένα τετράγωνο το οποίο αποτελεί τη ρίζα. Για την περικοπή του χώρου αναδρομικά χρησιμοποιείται το αντικείμενο (object) **Bounding Box**, το οποίο οριοθετεί τα σύνορα του τετραγώνου που ορίζει ο κάθε κόμβος περιέχοντάς τα ως **properties**.

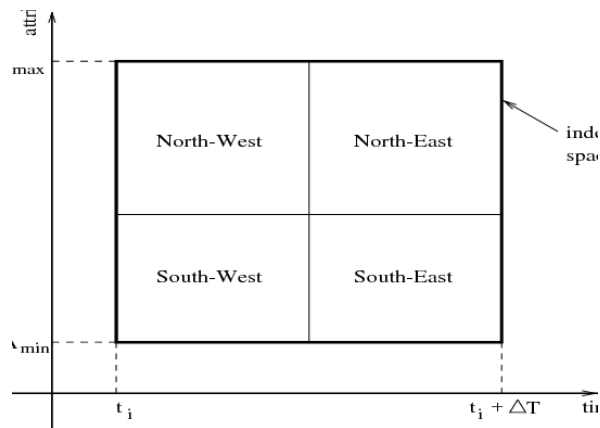
def contains():

Boolean συνάρτηση η οποία αληθεύει εφόσον το σημείο βρίσκεται μέσα στο bounding box που ελέγχθηκε.

Class QuadNode():

Οι μέθοδοι της κλάσης αυτής καθορίζουν τους κόμβους του δέντρου. Ως properties τίθενται το **κέντρο** του Quadtree, το **ύψος**, το **βάθος** του αλλά και την **χωρητικότητα** αυτού. Με την χωρητικότητα γίνεται έλεγχος του βασικού κανόνα του Quadtree, όπου ο αριθμός των παιδιών του να μην ξεπερνά τα 4. Για το ψάξιμο του δέντρου χρησιμοποιείται η κατεύθυνση του κόμβου. Έτσι εξετάζεται αν βρίσκεται εκεί που επιθυμείτε αναζήτηση.

Figure 3.



def contains(point):

Boolean συνάρτηση που επαληθεύεται στην περίπτωση που ένα σημείο (point) βρίσκεται μέσα στα δεδομένα του κόμβου.

def iter():

Μετατρέπει όλα τα σημεία (points) σε επαναληπτικά (iterable). Γι' αυτό χρησιμοποιείται η συνάρτηση **yield** της python. Η δήλωση yield επιστρέφει ένα αντικείμενο γεννήτριας (generator) σε αυτόν που καλεί τη συνάρτηση που περιέχει yield, αντί να επιστρέφει απλώς μια τιμή όπως κάνει η **return** κλήση.

def calc bounding box():

Υπολογισμός του τετραγώνου μέσα στο οποίο εμπεριέχεται ο κόμβος.

def contains_point(point):

Συνάρτηση τύπου Boolean που ελέγχει αν ένα σημείο (point) υπάρχει μέσα σε έναν κόμβο.

def is_ul(), def issuer(), def is_ll(), def is_lr():

Σε ποιον προσανατολισμό βρίσκεται το σημείο (point) ενδιαφέροντος (NE, NW, SE, SW)(Figure 3).

def subdivide():

Αναδρομική σχέση που υποδιαιρεί έναν κόμβο και τα παιδιά του.

def insert():

Boolean μέθοδος, που αληθεύει στην εισαγωγή του σημείου (point). Πρώτα γίνεται έλεγχος της χωρητικότητας του κόμβου, στον οποίο εισάγεται το σημείο (point). Αν το έχει υπερβεί τότε υποδιαιρείται. Ακολούθως εξετάζεται ο προσανατολισμός του και αναλόγως τοποθετείται.

def find():

Ψάχνει έναν κόμβο όπου το περιεχόμενο του σημείου (point) βρίσκεται στα παιδιά του. Επιστρέφει τον κόμβο και τα περιεχόμενά του.

def find_node():

Ψάχνει το κόμβο που μπορεί να περιέχει το σημείο (point). Ως όρισμα δέχεται τη λίστα των σημείων που έχει ήδη επισκεφτεί ο αλγόριθμος και το σημείο ενδιαφέροντος. Επιστρέφει ένα **QuadNode()** αντικείμενο (object) αν βρέθηκε και τη λίστα με όσα πέρασε για να το βρει.

def all_points():

Επιστρέφει μια λίστα με όλα τα σημεία (points) που περιέχει ένας κόμβος και τα παιδιά του.

def within_bb():

Αναδρομική συνάρτηση για τον έλεγχο ενός τετραγώνου αν περιέχεται μέσα σε κάποιο μεγαλύτερο, συμπεριλαμβανομένου και των πιθανών κατευθύνσεων που μπορούν να ακολουθήσουν.

def Euclidean distance():

Υπολογισμός της Ευκλείδειας απόστασης 2 σημείων στο χώρο.

Class QuadTree():

Δημιουργείται το αντικείμενο (object) **QuadTree**, όπου μέσα σε αυτό βρίσκονται κόμβοι και συντεταγμένες. Για την υλοποίηση του QuadTree απαιτείται η γνώση του κέντρου του, του ύψους του αλλά και του πλάτους του στο χώρο. Για τη χωρητικότητα η καθορισμένη (default) τιμή της ορίζεται σε **none**.

def convert to point():

Ελέγχεται η είσοδος (input) και επιστρέφεται με τον αποδεκτό τύπο (type).

def contains():

Boolean συνάρτηση που αληθεύει εφόσον ένα σημείο (point) περιέχει τη μέθοδο find().

def len():

Μήκος δέντρου.

def iter():

Μετατρέπει το δέντρο σε επαναληπτικό (iterable).

def insert():

Η είσοδος (input) μετατρέπεται σε object point και γίνεται μεταφόρτωση σε αυτήν τα χαρακτηριστικά. Στη συνέχεια, καλείται η **insert** (ορίστηκε παραπάνω).

def find():

Ψάχνει το σημείο (point) και το επιστρέφει ως αποτέλεσμα κλήσης της μεθόδου find().

Local QTree.py:

Class LocalQuadTree(object QUADTREE):

Αφορά σε μία κλάση προτύπου για το QuadTree έτσι ώστε να γίνει εφαρμογή των k κοντινότερων γειτόνων (K-Nearest-Neighbors).

def query():

Μέθοδος για την εύρεση γειτονικών σημείων συγκριτικά με το σημείο ενδιαφέροντος. Χρησιμοποιείται η συνάρτηση **deque()** ή αλλιώς «ουρά διπλού άκρου» της βιβλιοθήκης **collections**, η οποία είναι μια δομή δεδομένων που επιτρέπει την προσθήκη ή αφαίρεση στοιχείων από κάθε άκρο σε σταθερό χρόνο, $O(1)$. Εδώ, χρησιμοποιείται έτσι ώστε να υλοποιήσουμε μια στοίβα με τα σημεία που μπορεί να βρίσκονται «μέσα» στο σημείο ενδιαφέροντος. Εκεί, φτιάχνεται μία λίστα όπου από αριστερά μπαίνουν τα χαρακτηριστικά που απαρτίζουν τον κόμβο, ενώ από δεξιά κρατάει ένα ιστορικό με τους κόμβους που ελέγχθηκαν (Figure 4).

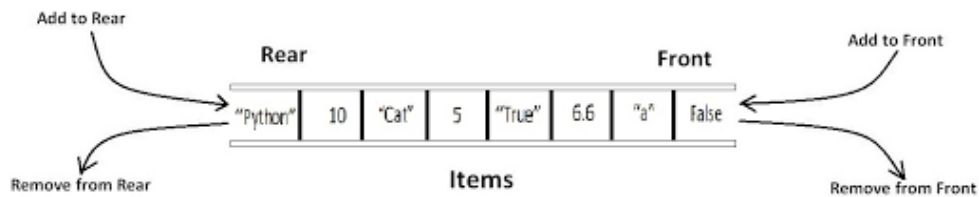


Figure 4

def get knn(point, k):

Ως attributes δέχεται το σημείο (point) της αναζήτησης μας και τον αριθμό των k κοντινότερων γειτόνων που απαιτείται να ληφθούν από τη στοίβα, την οποία επιστρέφει από την κλήση της μεθόδου **query()**. Εφαρμόζεται η Ευκλείδεια απόσταση και γίνεται ταξινόμηση των αποτελεσμάτων.

@staticmethod def nodeIsLeaf():

Μέθοδος που ανήκει στην κλάση **LocalQuadTree**, και όπως προοικονομεί το όνομά της, επιστρέφει αν ο κόμβος που εξετάζεται είναι φύλλο.

run QTree.py

Σε αυτό το αρχείο φορτώνεται το σύνολο των δεδομένων (dataset) και αποθηκεύονται κατάλληλα τα στοιχεία που απαιτούνται για την εφαρμογή των μεθόδων που υλοποιήθηκαν στα προηγούμενα εδάφια. Έπειτα, ο χρήστης μέσω ενός στοιχειώδους **menu** αποφασίζει αν θα εξετάσει τους k κοντινότερους γείτονες ενός επιστήμονα, ή το εύρος των κόμβων για κάποια συγκεκριμένα όρια του ορθογωνίου.

def put into list():

Αφού γίνει το διάβασμα του αρχείου, σε μία for λούπα τοποθετούνται τα στοιχεία του csv που μας ενδιαφέρουν στις λίστες latitude, longitude, scientist_list.

def get points(), range search():

Οι συγκεκριμένες συναρτήσεις παραμένουν ίδιες όσον αφορά την υλοποίησή τους, όπως και στις προηγούμενες πολυδιάστατες δομές που εξετάσαμε.

def run scientists():

Δέχεται ως είσοδο το όνομα του επιστήμονα και τα στοιχεία του συνόλου δεδομένων (dataset). Με σκοπό τη δημιουργία ενός **object point** για την εύρεση, λαμβάνει μέρος μία προ-επεξεργασία στη μέθοδο main_scientists(), στοχεύοντας στην εύρεση του **index** μέσα στη λίστα δεδομένων. Δίνεται η επιλογή για k-NN ή range_search. Στο σημείο αυτό, τονίζεται πως ο χρόνος ξεκινάει να μετράει με την έναρξη του k-NN και όχι με τη δημιουργία του δέντρου.

names.py:

def name filter(line, start, end):

Το συγκεκριμένο script δημιουργήθηκε για να φιλτράρονται τα ονόματα των επιστημόνων σύμφωνα με το πρώτο γράμμα που είναι καταχωρημένοι, σε ένα εύρος γραμμάτων (start, end).

Αποτελέσματα:

- K Nearest Neighbors:

```
Reading scientists data...

Give the 1st letter of list: G
Give the last letter of list: Z
Data successfully read.

MENU :

1: kNN
2: Range Search
1
Give the name of the scientist you want to search his neighbors: Victor Bahl
How many nearby scientists do you want to show: 5
*****
For Victor Bahl with coordinates 462, 4

*****
The algorithm run for :
0.002977sec
*****
*****
The 5 nearest neighbors are :

0 ==> Gerard Holzmannwith coordinates (463,4)
1 ==> Yanhong Annie Liuwith coordinates (461,6)
2 ==> Samir Daswith coordinates (464,6)
3 ==> Wendy Hallwith coordinates (459,3)
4 ==> Matthew Dillonwith coordinates (460,1)
```

- Range Search:

Reading scientists data...

Give the 1st letter of list: G

Give the last letter of list: Z

Data successfully read.

MENU :

1: kNN

2: Range Search

2

Please give the minimum cord: 50

Please give the minimum award: 4

Please give the maximum cord: 150

Please give the maximum award: 6

Samson Abramsky[104, 5]

Sanjeev Arora[70, 4]

Rod Canion[87, 5]

John Carmack[51, 5]

Shimon Even[97, 5]

Raphael Finkel[86, 5]

Robert Floyd[134, 6]

Robert France[111, 4]

Kunihiko Fukushima[144, 5]

Robert M. Graham[123, 4]

Shelia Guberman[50, 4]

Jo rg Gutknecht[93, 4]

Margaret Hamilton[72, 4]

Johan Ho stad[117, 4]

Kenneth E. Iverson[95, 5]

William Kahan[56, 6]

Ken Kennedy[143, 5]

Leonard Kleinrock[101, 5]

Leslie Lamport[65, 6]

Nancy Lynch[115, 5]

Nadia Magnenat Thalmann[148, 6]

John McCarthy[75, 5]

Marshall Kirk McKusick[71, 5]

Robin Milner[133, 6]

J Strother Moore[84, 5]

Peter Naur[59, 5]

Paritosh Pandya[106, 6]

Juan Pavo n[100, 4]

Jon Postel[142, 4]

Maciej Stachowiak[96, 5]

Richard Stallman[119, 6]

Richard E. Stearns[78, 4]

Robert Tarjan[124, 6]

Paul Vixie[76, 5]

Kevin Warwick[54, 4]

Joseph Weizenbaum[80, 4]

John Yen[127, 5]

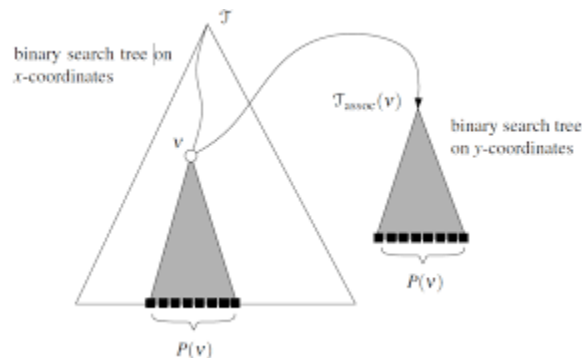
Shlomo Zilberstein[125, 6]

Range Search Tree

Γενική εισαγωγή:

Ένα Range Search Tree είναι μία δομή δεδομένων που χρησιμοποιείται για την εύρεση και την καταγραφή ενός συνόλου από σημεία σε ένα πολυδιάστατο χώρο, για ένα δεδομένο εύρος. Γενικότερα ένα range tree d-διαστάσεων αποτελεί ένα αναδρομικά καθορισμένο πολυεπίπεδο δυαδικό δέντρο αναζήτησης. Αυτό σημαίνει πως κάθε level της δομής, αναπαρίσταται από ένα BST σε μία από τις διαστάσεις d. Κατά τη διάρκεια της αναζήτησης, διασχίζεται το δέντρο ώστε να βρεθούν όλοι οι επικαλυπτόμενοι κόμβοι μιας δεδομένης περιοχής και να επιστραφούν σαν αποτέλεσμα τα points που είναι αποθηκευμένα σε κάθε κόμβο ξεχωριστά.

Επιπλέον ο χρόνος εκτέλεσης είναι λογαριθμικός $O(\log^d n + k)$, καθιστώντας την μια αρκετά αποτελεσματική δομή, συγκριτικά και με τα KD-trees, ειδικά όταν παρουσιαστούν μεγάλος όγκος δεδομένων.



range tree.py:

Στο αρχείο γίνεται η κατασκευή ολόκληρου του δέντρου με την κλήση της range search, καθώς και μερικών βοηθητικών της συναρτήσεων.

Αρχικά ορίζουμε μία μεταβλητή με όνομα DIMENSIONS στην οποία αποθηκεύουμε τις διαστάσεις (2) του δέντρου.

Class Node :

Δημιουργία αντικειμένου Node, για κάθε κόμβο, με τα παρακάτω properties:

1. coords
2. name
3. left
4. right
5. next_dimension

Create Range Tree:

Η συγκεκριμένη συνάρτηση συμβάλλει στην κατασκευή του δέντρου. Λαμβάνει ως είσοδο μια λίστα διατεταγμένων συντεταγμένων των κόμβων και επιστρέφει πίσω σαν έξοδο τη ρίζα του δέντρου, καθώς και μια νέα λίστα των coords για τη κατασκευή του δέντρου της επόμενης διάστασης (αν υπάρχει).

Η διαδικασία της αποτύπωσης του δέντρου ξεκινά από τη λίστα με τις διατεταγμένες συντεταγμένες και επιλέγοντας ως ρίζα του δέντρου το στοιχείο που χωρίζει τη λίστα στα δύο. Με αυτό τον τρόπο δημιουργείται ένα αριστερό και ένα δεξί υποδέντρο αντίστοιχα, όπου για κάθε ένα ξεχωριστά καλούμε αναδρομικά τη συνάρτηση. Χρησιμοποιούμε την (αμέσως επόμενη) συνάρτηση merge, ώστε να συνθέσουμε τις δυο λίστες μαζί με την εκάστοτε ρίζα, ως προς την επόμενη διάσταση και να δημιουργήσουμε το δέντρο που αντιστοιχεί σε αυτήν.

Merge:

Η συνάρτηση merge, που παίρνει ως ορίσματα τη λίστα του αριστερού και του δεξιού υποδέντρου και την ρίζα για την επόμενη διάσταση. Προκειμένου να βρεθεί η σωστή θέση της ρίζας στην ενιαία λίστα, γίνεται πρώτα στοίχιση των δύο λιστών σύμφωνα με τους δείκτες.

Search:

Με τη συγκεκριμένη συνάρτηση βρίσκουμε ένα συγκεκριμένο κόμβο. Το traversal πραγματοποιείται συγκρίνοντας τις συντεταγμένες του κόμβου που ψάχνουμε με της ρίζας του δέντρου.

Find split node:

Είναι η βοηθητική συνάρτηση για την range search. Εύρεση του κόμβου split_node με αναδρομική κλήση της συνάρτησης ελέγχοντας κάθε φορά τις συντεταγμένες της ρίζας με τα min, max από το search path.

Is in range:

Ως είσοδο δέχεται συντεταγμένες κόμβου και το διάστημα αναζήτησης, ενώ στην έξοδο επιστρέφει Boolean τιμές ανάλογα με το αν οι δοθείσες συντεταγμένες ανήκουν στο διάστημα.

Range Search:

Με τη βοήθεια των συναρτήσεων Find_split_node, Is_in_range πραγματοποιούμε δύο δυαδικές αναζητήσεις για κάθε διάσταση. Κατευθυνόμενοι προς τα αριστερά, και άρα προσεγγίζοντας τη min τιμή του διαστήματος, θα πρέπει παράλληλα να κάνουμε store στη λίστα τα points που βρίσκονται στα δεξιά υποδέντρα και αντίστοιχα, όσο προσεγγίζουμε τη max τιμή να αποθηκεύουμε τους κόμβους που βρίσκονται στα αριστερά υποδέντρα από αυτά της αναζήτησης.

menu.py:

Για την εκτέλεση του κώδικα στο αρχείο menu.py, κάνουμε import το αρχείο range_tree.py ώστε να φτιαχτεί το τελικό αρχείο που θα τρέξουμε για να πάρουμε τα αποτελέσματα. Πρώτα όμως διαβάζουμε το αρχείο csv, παραλείποντας την πρώτη γραμμή που βρίσκεται η επικεφαλίδα, και παράλληλα δημιουργούμε κόμβους ώστε να εφαρμόσουμε την αναζήτηση.

Για τα πειράματά μας, έχουν φτιαχτεί 4 επιπλέον συναρτήσεις:

1. Pre order:

Διασχίζουμε τους κόμβους του δέντρου από πάνω προς τα κάτω και από αριστερά προς τα δεξιά.

2. Print nodes:

Εκτυπώνει τους κόμβους του return της range_search.

3. Euclidean distance:

Η ευκλείδεια απόσταση μεταξύ δύο κόμβων για την υλοποίηση του kNN από τον τύπο $d(x,y)=\sqrt{\sum (y_i-x_i)^2}$

4. kNN algorithm:

Στόχος του αλγορίθμου είναι να βρεθούν οι γείτονες ενός κόμβου στο δέντρο.

Σύμφωνα με τους κανόνες της δομής, οι κόμβοι του δέντρου είναι διατεταγμένοι, που σημαίνει πώς γύρω του θα βρίσκονται κόμβοι με κοντινές συντεταγμένες.

Επομένως, αρχικοποιούμε τυχαία ένα range, καλώντας την συνάρτηση range_search, ώστε να αντλήσουμε τους κόμβους που βρίσκονται στην άμεση γειτονιά του κόμβου. Τέλος υπολογίζουμε την απόσταση, σύμφωνα με την euclidean_distance και ταξινομούμε τους κόμβους, και αν οι γειτονικοί κόμβοι είναι λιγότεροι από k, επεκτείνουμε το διάστημα αναζήτησης και επαναλαμβάνουμε τη διαδικασία.

Αποτελέσματα:

Αρχικά δημιουργείται το δέντρο με σύνολο κόμβων **670**, ένα για κάθε επιστήμονα, και αμέσως μετά εμφανίζεται το περιβάλλον του μενού για το χρήστη με τις παρακάτω δυνατές επιλογές:

- Print tree (0)
- Range Search (1)
- kNN (2)
- Exit program (-1)

Ανάλογα με τον αριθμό που θα πληκτρολογήσει ο user, το πρόγραμμα εκτελεί την αντίστοιχη λειτουργία.

Παρακάτω επισυνάπτονται τα στιγμιότυπα μετά από την εκτέλεση του αρχείου menu.py:

Range – Search:

```
1
Give minimum coordinate for dimension 0
10
Give maximum coordinate for dimension 0
100
Give minimum coordinate for dimension 1
2
Give maximum coordinate for dimension 1
4
-----
```

Nodes found (37)

[83.0, 2.0]		Name:Peter O'Hearn
[41.0, 3.0]		Name:David P. Reed
[72.0, 4.0]		Name:Margaret Hamilton
[54.0, 4.0]		Name:Kevin Warwick
[43.0, 4.0]		Name:Jochen Liedtke
[50.0, 4.0]		Name:Shelia Guberman
[67.0, 4.0]		Name:Fernando J. Corbat
[70.0, 4.0]		Name:Sanjeev Arora
[55.0, 2.0]		Name:Atta ur Rehman Khan
[48.0, 3.0]		Name:Aaron Sloman
[80.0, 4.0]		Name:Joseph Weizenbaum
[78.0, 4.0]		Name:Richard E. Stearns
[62.0, 3.0]		Name:Sergey Brin
[52.0, 3.0]		Name:Philip Woodward
[74.0, 3.0]		Name:Roberto Tamassia
[66.0, 3.0]		Name:James D. Foley
[82.0, 3.0]		Name:Jaime Teevan
[57.0, 2.0]		Name:Robert Kowalski
[68.0, 2.0]		Name:Grady Booch
[60.0, 2.0]		Name:Martin Odersky
[20.0, 2.0]		Name:Mary Allen Wilkes
[37.0, 3.0]		Name:Jon Crowcroft
[36.0, 4.0]		Name:Dennis Ritchie
[29.0, 4.0]		Name:Ivar Jacobson
[40.0, 4.0]		Name:Beatrice Helen Worsley
[38.0, 3.0]		Name:Dan Klein
[33.0, 2.0]		Name:Seinosuke Toda
[10.0, 3.0]		Name:David F. Bacon
[15.0, 4.0]		Name:Stephen Wolfram
[11.0, 2.0]		Name:David Patterson
[12.0, 4.0]		Name:John Krogstie
[91.0, 3.0]		Name:Alston Householder
[93.0, 4.0]		Name:Jörg Gutknecht
[85.0, 3.0]		Name:John O'Sullivan
[99.0, 2.0]		Name:Edward Yourdon
[98.0, 2.0]		Name:David J. Brown
[100.0, 4.0]		Name:Juan Pavón

kNN algorithm:

```
Give your Scientist's name: Rediet Abebe
Number of scientists you want: 5
-----
|Distance:0.0 |   Name:Rediet Abebe   |   Coordinates:[325.0, 5.0]

|Distance:4.0 |   Name:Bernard Richards   |   Coordinates:[323.0, 5.0]

|Distance:5.0 |   Name:Stephen Cook   |   Coordinates:[326.0, 3.0]

|Distance:8.0 |   Name:Gene Amdahl   |   Coordinates:[327.0, 3.0]

|Distance:10.0 |   Name:Harry Bouwman   |   Coordinates:[328.0, 4.0]

|Distance:10.0 |   Name:Jeremy Gibbons   |   Coordinates:[324.0, 2.0]
```

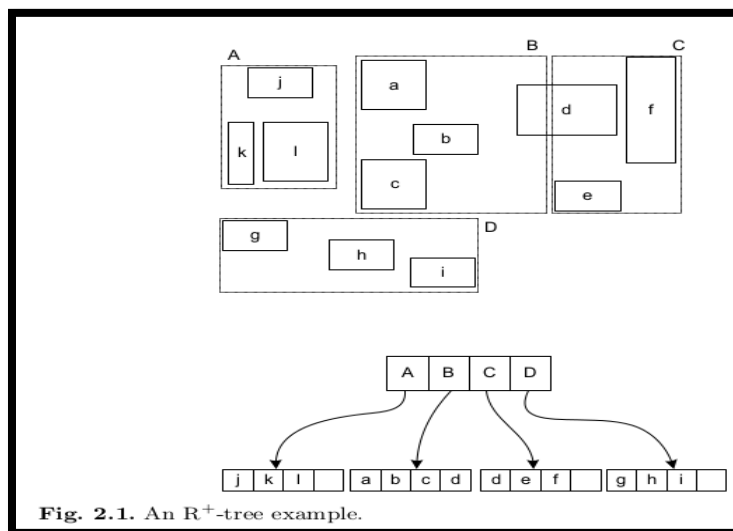
***Technical report:** Δεν πραγματοποιείται σωστή αποτύπωση του δέντρου στην επιλογή Print tree και για το λόγο αυτό δεν συμπεριλήφθηκε στην αναφορά.

R – Tree

Εισαγωγή:

Ένα R-tree είναι μια δομή δεδομένων που χρησιμοποιείται για χωρική ευρετηρίαση πολυδιάστατων δεδομένων, που συνήθως χρησιμοποιείται για την οργάνωση χωρικών δεδομένων σε δύο ή περισσότερες διαστάσεις. Σε ένα δέντρο R, κάθε κόμβος αντιπροσωπεύει ένα ορθογώνιο που περιέχει ένα ή περισσότερα αντικείμενα. Το R-δέντρο διαιρεί αναδρομικά το χώρο σε μικρότερες ορθογώνιες περιοχές έως ότου κάθε ορθογώνιο περιέχει μικρό αριθμό αντικειμένων ή φτάσει σε ένα ελάχιστο μέγεθος. Κάθε κόμβος στο δέντρο είναι είτε ένας κόμβος φύλλου, ο οποίος περιέχει τα πραγματικά αντικείμενα, είτε ένας εσωτερικός κόμβος, ο οποίος περιέχει δείκτες σε άλλους κόμβους.

Τα R-trees χρησιμοποιούνται συνήθως για χωρική ευρετηρίαση σε συστήματα γεωγραφικών πληροφοριών και άλλες εφαρμογές που ασχολούνται με μεγάλο όγκο χωρικών δεδομένων, όπως υπηρεσίες που βασίζονται σε τοποθεσία, εξόρυξη δεδομένων και όραση υπολογιστή.



Rtree.py:

Στο αρχείο rtree.py βρίσκονται όλες οι συναρτήσεις για τη δημιουργία του δέντρου, αλλά και για τις απαραίτητες πράξεις:

length(data, j, D):

Η μέθοδος αυτή παίρνει σαν είσοδο την λίστα με τις συντεταγμένες τα minnodes που μας δείχνει στα πόσα τμήματα θα μοιραστεί η λίστα μας και το πλήθος των κόμβων που περιέχεται στη λίστα.

min_max(List, data):

Η μέθοδος αυτή παίρνει σαν είσοδο την λίστα με τις συντεταγμένες και τη λίστα με τα τμήματα που θα μοιραστεί η λίστα. Αυτή η μέθοδος υπολογίζει τα νέα bounding box που θα προκύψουν από τα νέα minx , miny , xmax , ymax και τα προσθέτει στο δέντρο.

make_tree(self, data):

Η μέθοδος αυτή καλείται για να δημιουργηθεί το δέντρο. Η μέθοδος αυτή ελέγχει σε ποιο υπό-bounding box ανήκει το νέο bounding box.

dist(self, node):

Η μέθοδος αυτή καλείται όταν χρειάζεται να υπολογιστεί η απόσταση από ένα σημείο.

clear_temp(self, k):

Η μέθοδος αυτή καλείται μετά από κάθε εισαγωγή στο δέντρο για να καθαρίσει το πίνακα temp από τα δεδομένα που δεν χρειάζονται πλέον.

find_node(self, node):

Η μέθοδος αυτή καλείται όταν ο χρήστης επιλέξει να βρει αν υπάρχει κάποιο σημείο στο δέντρο. Στη περίπτωση που υπάρχει επιστρέφει Node found.

range_Search(self,node):

Η μέθοδος αυτή καλείται όταν ο χρήστης επιλέξει να δει πόσοι επιστήμονες υπάρχουν σε εκείνο το Bounding box.

Αποτελέσματα:

Main Menu:

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
```

1. Find node :

Ζητείται από το χρήστη να βάλει τον αριθμό των cords και των awards . Αμέσως ψάχνει αν υπάρχει αυτό το σημείο. Σε περίπτωση που υπάρχει τυπώνει **Node found** αλλιώς τυπώνει **Node didn't exist**.

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
1
Give node coordinates you want to find:
Give cord first:
600
Then award:
1
Judea Pearl
```

2. Get n nearest :

Ζητείται από το χρήστη να βάλει τον αριθμό cords , τον αριθμό awards και τον αριθμό των πλησιέστερων γειτόνων που θέλει να βρεθούν.

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
2
Give node coordinates you want to find:
Give cord first:
150
Then award :
5
Give number of nearest:
2
Peter Norvig
Kristen Nygaard
.....
```

3. Range search

Ζητείται από το χρήστη να εισάγει ελάχιστο και τον μέγιστο αριθμό awards και cords που επιθυμεί για να δημιουργηθεί το Bounding box και να βρει πόσοι επιστήμονες είναι μέσα στο Bounding box.

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
3
Give cordmin
111
Give awardmin:
2
Give cordmax:
130
Give awardmax:
3
Jeff Bonwick
John Romero
Christopher Strachey
Edward H. Shortliffe
Jack Dongarra
Klara Dan von Neumann
Jan Dietz
```

4. Search by [A-Z], awards, and cords:

Ζητείται από το χρήστη να βάλει τα αρχικά των ονομάτων για τους επιστήμονες, μαζί με τα awards και cords, έτσι ώστε να τυπωθούν οι επιστήμονες που πληρούν όλα τα κριτήρια.

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
4
Enter the First character: K
Enter the Second character: M
Enter the max award:4
Enter the max cord:150
The names in the namelist with awards less than or equal to 4 and with cords less than or equal to 150 sorted by character:
Give choice what you want to do:

, 'Martin Odersky ', 'Mary Allen Wilkes ', 'Munindar P. Singh ', 'Murray Turoff ']
```

```
: ['Kevin Warwick ', 'Klara Dan von Neumann', 'Kristen Nygaard ', 'Luca Cardelli ', 'Luis von Ahn ', 'Manindra Agrawal ', 'Margaret Burnett ', 'Margaret Hamilton ',
```

5. ADD Node:

Ζητείται από το χρήστη να βάλει τα awards και cords, αλλά και το όνομα που επιθυμεί να προσθέσει.

```
Give choice what you want to do:
1:Find node
2:Get n nearest
3:Range search
4:Search by [A-Z] , awards and cords
5:Add node
6:Exit
5
Give cords:
500
Then awards:
50
And Name:
Marios Priamos
[500, 50, 'Marios Priamos']
```

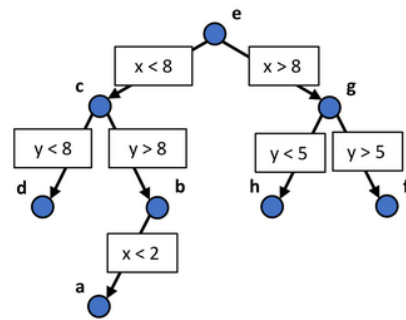
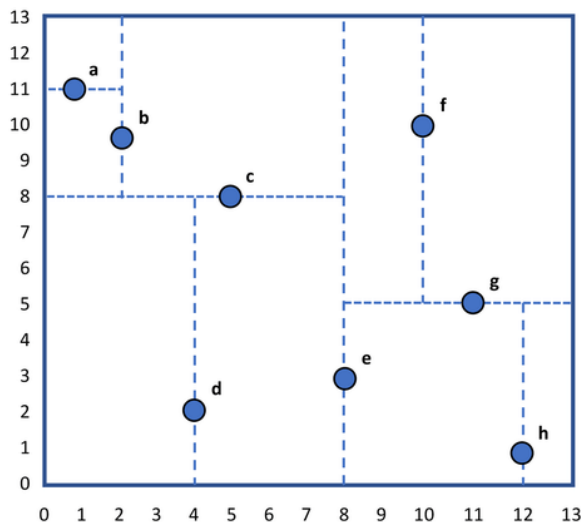

K-D Tree:

Εισαγωγή:

Ένα δέντρο k-d είναι μια δομή δεδομένων που χρησιμοποιείται συνήθως για αναζήτηση πλησιέστερου γείτονα σε χώρους υψηλών διαστάσεων. Το k-d σημαίνει k-dimensional, το οποίο αναφέρεται στο γεγονός ότι τα σημεία δεδομένων αναπαρίστανται ως σημεία σε ένα k-διάστατο χώρο.

Το δέντρο k-d είναι ένα δυαδικό δέντρο, όπου κάθε κόμβος αντιπροσωπεύει ολόκληρο τον χώρο και κάθε επόμενο επίπεδο του δέντρου χωρίζει το χώρο σε δύο μισά κατά μήκος μιας από τις διαστάσεις k. Η κατάτμηση επιτυγχάνεται με την επιλογή ενός σημείου περιστροφής, το οποίο χρησιμοποιείται για να χωρίσει το σμήγμα σε δύο υποχώρους.

Η διαδικασία κατασκευής ενός δέντρου k-d περιλαμβάνει την αναδρομική κατάτμηση του χώρου έως ότου κάθε κόμβος αντιπροσωπεύει ένα μόνο σημείο δεδομένων ή έναν μικρό αριθμό σημείων δεδομένων. Μόλις κατασκευαστεί το δέντρο, μπορεί να χρησιμοποιηθεί για την αποτελεσματική αναζήτηση του πλησιέστερου γείτονα ενός δεδομένου σημείου ερωτήματος στον k-διάστατο χώρο.



Kd-tree.py:

1. function kdtree:

Εδώ δημιουργείται το k-d δέντρο. Αν η λίστα με τα points είναι άδεια τότε επιστρέφει τη ρίζα, ενώ αν δεν είναι άδεια τότε κατασκευάζει το k-d δέντρο.

2. function insert:

Πρώτα μετράμε την απόσταση μεταξύ του κάθε κόμβου που υπάρχει στο δέντρο με τον καινούργιο κόμβο. Αν η απόσταση είναι μικρότερη του 0 τότε τοποθετείται ο καινούργιος κόμβος και αν είναι μεγαλύτερη ή ίση του 0 ξανακαλείται η συνάρτηση.

3. function delete:

Ελέγχουμε αν υπάρχει το point στη λίστα, αν υπάρχει τότε διαγράφεται από τη λίστα και ξαναδημιουργούμε το δέντρο. Αν όχι επιστρέφει το κατάλληλο μήνυμα.

4. function getknn:

Εδώ δημιουργούμε ένα heap όπου εισάγουμε την απόσταση μεταξύ ενός κόμβου και των γειτόνων του. Αν το μέγεθος του heap είναι μικρότερο από των αριθμών των γειτόνων που αναζητούμε τότε εισάγουμε στο σωρό το κόμβο-γείτονα και την απόσταση του. Αν η απόσταση είναι μικρότερη του κόμβου στο τελευταίο σημείο του heap, τότε αυτό το σημείο το αφαιρούμε. Αυτό επαναλαμβάνεται μέχρι το μέγεθος του heap να ισούται με τον αριθμό των γειτόνων που αναζητούμε και επιστρέφει τους γείτονες.

5. function range search:

Παίρνει ως ορίσματα το διάστημα που θέλουμε να ψάξουμε(range_min, range_max) και μέσω αυτής γίνεται έλεγχος μέσω της λίστας και επιστρέφει το αποτέλεσμα.

6. function main:

Εδώ καλούμε τις προηγούμενες συναρτήσεις που αναφέραμε. Χρησιμοποιούμε την βιβλιοθήκη pandas για να διαβάσουμε το demo_file.csv και δημιουργούμε τη λίστα με τα points και ακόμη μέσω του pd.transform κάνουμε pre-process τη στήλη education για να έχουμε integers αντί strings. Επίσης χρησιμοποιούμε το pprint για να έχουμε μια εικόνα για το δέντρο.

Αποτελέσματα:

- Κατασκευάζουμε το δέντρο με την εντολή `tree = kdtree(points)` και το εμφανίζουμε με το `pp.print(tree)`:

***Εδώ βάζουμε ένα μέρος του αποτελέσματος.**

```
E:\anaconda3\envs\tf\python.exe C:\Users\fear1\Desktop\polidiastates\project2022\kd-tree.py
[ [ [ [ [ [ [ [ [None, None, [1, 1]],
      [None, None, [6, 1]],
      [3, 1]],
      [[None, None, [10, 1]], None, [12, 1]],
      [7, 1]],
    [ [[None, None, [19, 1]], None, [21, 1]],
      [[None, None, [29, 1]], None, [30, 1]],
      [24, 1]],
      [16, 1]],
```

- Το `insert(tree, test_point)` προσθέτει το σημείο `[0, 1]` στο δέντρο μας και μετά καλούμε το `pp.print(tree)`.

```
[ [ [ [ [ [ [ [ [ [None, None, [0, 1]], None, [1, 1]],
      [None, None, [6, 1]],
      [3, 1]],
      [[None, None, [10, 1]], None, [12, 1]],
      [7, 1]],
    [ [[None, None, [19, 1]], None, [21, 1]],
      [[None, None, [29, 1]], None, [30, 1]],
      [24, 1]],
      [16, 1]],
```

- Το `pp.print(delete(points, test_points))` αφαιρεί το σημείο `[1, 1]` από την λίστα και εκτυπώνει το δέντρο.

```
E:\anaconda3\envs\tf\python.exe C:\Users\fearl\Desktop\polidiastates\project2022\kd-tree.py
[ [ [ [ [ [ [ [ [None, None, [3, 1]],
      [None, None, [7, 1]],
      [5, 1]],
      [[None, None, [12, 1]], None, [13, 1]],
      [8, 1]],
      [ [None, None, [21, 1]], None, [23, 1]],
        [None, None, [30, 1]], None, [31, 1]],
        [28, 1]],
        [17, 1]],
      [1, 1]],
    [1, 1]],
  [1, 1]],
[1, 1]]
```

- Το `print(get_knn(tree, [6, 1], 4))` μας εμφανίζει 4 γείτονες του σημείου:

```
E:\anaconda3\envs\tf\python.exe C:\Users\fearl\Desktop\polidiastates\project2022\kd-tree.py
[(0, [6, 1]), (1, [7, 1]), (3, [3, 1]), (4, [10, 1])]
```

- Τέλος το `print(range_search(points, [1,1], [6,1]))` μας εμφανίζει τα σημεία που βρίσκονται μεταξύ των δύο σημείων:

```
E:\anaconda3\envs\tf\python.exe C:\Users\fearl\Desktop\polidiastates\project2022\kd-tree.py
[[5, 1], [1, 1], [6, 1], [3, 1], [4, 1], [2, 1]]
```

- Τέλος -