



ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ

Εθνικόν και Καποδιστριακόν
Πανεπιστήμιον Αθηνών

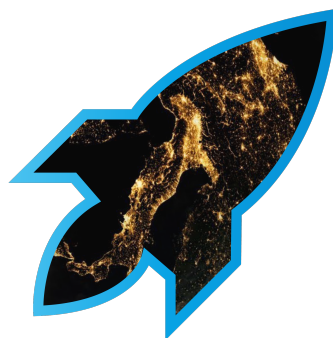
— ΙΔΡΥΘΕΝ ΤΟ 1837 —

ΑΝΑΠΤΥΞΗ ΛΟΓΙΣΜΙΚΟΥ ΓΙΑ ΑΛΓΟΡΙΘΜΙΚΑ ΠΡΟΒΛΗΜΑΤΑ

ΧΕΙΜΕΡΙΝΟ ΕΞΑΜΗΝΟ 2020

1η ΠΡΟΓΡΑΜΜΑΤΙΣΤΙΚΗ ΕΡΓΑΣΙΑ

ΑΝΑΖΗΤΗΣΗ ΚΑΙ ΣΥΣΤΑΔΟΠΟΙΗΣΗ ΔΙΑΝΥΣΜΑΤΩΝ ΣΤΗ C/C++



Αριθμός Μητρώου(ΑΜ):

1115201700217

1115201700203

Ονοματεπώνυμο:

Ορέστης ΣΤΕΦΑΝΟΥ

Λεωνίδας ΕΦΡΑΙΜ

ΑΚΑΔΗΜΑΪΚΗ ΧΡΟΝΙΑ 2020-2021

ΠΕΡΙΕΧΟΜΕΝΑ

1	ΕΙΣΑΓΩΓΗ	3
2	ΜΕΤΑΓΛΩΤΤΙΣΗ-ΕΚΤΕΛΕΣΗ	4
3	ΥΛΟΠΟΙΗΣΗ	5
3.1	ΕΙΣΟΔΟΣ ΔΕΔΟΜΕΝΩΝ	5
3.2	ΜΕΤΡΙΚΕΣ	6
3.3	HASH TABLE	6
3.4	LSH	7
3.5	HYPER CUBE	8
3.6	CLUSTERING	10
3.6.1	Lloyds	10
3.6.2	LSH Range Search	10

ΕΙΣΑΓΩΓΗ

Στα πλέσια της εργασίας είχαμε να υλοποιήσουμε τον αλγόριθμο LSH για διανύσματα στον D-διάστατο χώρο, καθώς και τον αλγόριθμο τυχαίας προβολής στον υπερκύβο βάσης της μετρικής Μανχάταν L1. Στην συνέχεια έπρεπε να εκτελέσουμε κάποια queries στο dataset που μας δώθηκε έτσι ώστε να επαληθεύσουμε την σωστή λειτουργία των αλγορίθμων. Τέλος κληθήκαμε να υλοποιήσουμε τους αλγόριθμους για την συσταδοποίηση διανυσμάτων βάση της μετρικής Μανχάταν όπου η ανάθεση θα έπρεπε να γίνει με τον αλγόριθμο του Lloyd's ή με αντίστροφη ανάθεση μέσω Range Search με LSH. Η υλοποίηση της εργασίας έχει γίνει σε C++

ΜΕΤΑΓΛΩΤΤΙΣΗ-ΕΚΤΕΛΕΣΗ

Για τις ανάγκες της εργασίας δημιουργήσαμε 3ις main συναρτήσεις όπου οι δύο είναι υπεύθυνες για του αλγόριθμους LSH και Hypercube, ενώ η τρίτη είναι υπεύθυνη για το Clustering

Η μεταγλώττιση γίνεται με τις παρακάτω εντολές

- **make lsh**
- **make cube**
- **make cluster**

Ενώ η εκτέλεση των προγραμμάτων γίνεται με τις εντολές που μας δώθηκαν στην εκφώνηση της εργασίας, δηλαδή:

- **LSH**

```
./lsh -d <input file> -q <query file> -k <int> -L <int> -o <output file> -  
N<number of nearest> -R <radius>
```

- **HYPER CUBE**

```
./cube -d <input file> -q <query file> -k <int> -M <int> -probes <int>  
-o<output file> -N <number of nearest> -R <radius>
```

- **CLUSTERING**

```
./cluster -i <input file> -c <configuration file> -o <output file> -complete  
<optional> -m <method: Classic OR LSH or Hypercube>
```

ΥΛΟΠΟΙΗΣΗ

3.1 ΕΙΣΟΔΟΣ ΔΕΔΟΜΕΝΩΝ

Για την εισαγωγή των δεδομένων έχουμε δημιουργήσει μια συνάρτηση με το όνομα **ReadData** η οποία δέχεται σαν όρισμα το path με το αρχείο εικόνων και ένα vector όπου στην συνέχεια το γεμίζει με τις εικόνες.

Η συνάρτηση αφού ανοίξει το αρχείο διαβάζει διαδοχικά 4 integers όπου αντιπροσωπεύουν αντιστοιχία

- Το magic number
- Το ύψος της εικόνας της εικόνας
- Το πλάτος της εικόνας της εικόνας
- Τον αριθμό των εικόνων που υπάρχουν στο αρχείο

Αφού ξέρουμε τις διαστάσεις των εικόνων τώρα μπορούμε να διαβάζουμε $N*N$ chars και να τους αποθηκεύουμε σε μια γραμμή του vector διαδοχικά.

Για την υλοποίηση της συνάρτησης ReadData χρειαστήκαμε να υλοποιήσουμε ακόμα μια συνάρτηση με όνομα **NumReverse** η οποία πέρνει ένα integer και του αλλάζει το endian του με μερικά shifts γιατί ο αριθμός που υπάρχει στο αρχείο είναι ανάποδα οπότε πρέπει να αντιστραφεί.

3.2 ΜΕΤΡΙΚΕΣ

Για τις μετρικές δημιουργίσαμε μια κλάση με το όνομα **Metrics** η οποία έχει μια συνάρτηση με το όνομα **get_distance** η οποία δέχεται σαν όρισμα τις 2 εικόνες που θέλουμε να βρούμε της απόστασή τους καθώς και ακόμα ένα όρισμα το οποίο είναι το όνομα της μετρικής π.χ. L1 για την μετρική Μαχνάταν. Έδω μπορούν να υλοποιηθούν και άλλες μετρικές αλλά στην εργασία μας ζητήθηκαν μόνο η μετρική Μαχνάταν.

Για την υλοποίηση της Μαχνάταν μετρικής πήραμε το άθροισμα της απόλυτης τιμής των σημείων των δύο εικόνων

3.3 HASH TABLE

Για την υλοποίηση του Lsh χρειαστικάμε ένα hashtable οπότε δημιουργήσαμε μια κλάση με το όνομα hashtable. Η κλάση αυτή αποτελείται από τα buckets που είναι ένας πίνακας με vectors, το μέγεθος του πίνακα, τις σταθερές K και W, μια μεταβλητή sRandInit η οποία αρχικοποιεί την rand για να έχουμε τυχαία s κάθε φορά καθώς και ένα vector με vectors το οποίο περιέχει τα s. Στον **constructor** της hashtable αρχικοποιούμε όλες τις μεταβλητές καθώς και δημιουργούμε τα τυχαία s όπου τα βάζουμε στο vector. Εκτός από τον constructor το hashtable έχει και τις παράτακτω συναρτήσεις

- **hash_function**

Η συνάρτηση αυτή δημιουργεί την συνάρτηση $g(p)$ σύμφωνα με τον αλγόριθμο lsh. Αυτό το κάνει φτιάχνοντας μια διαφορετική $h(p)$ κάθε φορά βάσει του παρακάτω τύπου

$$h(p) = a_{d-1} + m \cdot a_{d-2} + \dots + m^{d-1} \cdot a_0 \bmod M \in \mathbb{N},$$

Στην συνέχεια ενώνει όλες τις $h(p)$ για να δημιουργήσει την $g(p)$

$$g(p) = [h_1(p)|h_2(p)|\dots|h_k(p)] \in \mathbb{N}.$$

- **insert**

Η συνάρτηση αυτή δέχεται σαν όριμα μια εικόνα καθώς και τον αριθμό του bucket που πρέπει να μπει με στόχο να εισάγει την εικόνα αυτή στο κατάλληλο bucket του hashtable

- **get_bucket_imgs**

Η συνάρτηση αυτή πέρνει σαν όρισμα τον αριθμό κάποιου bucket και επιστρέφει ένα vector με τα στοιχεία αυτού του bucket

3.4 LSH

Ο αλγόριθμος LSH υλοποιήτε μέσω μιας κλάσης με το αντίστοιχο όνομα. Η κλάση αυτή περιέχει τις σταθερές K, L, r , ένα hash table και ένα vector με όλα τα δεδομένα των εικόνων. Στον constructor αρχικοποιούνται όλες οι μεταβλητές. Επίσης δημιουργούνται όλα τα hashfunction και μπένει η κάθε εικόνα στο bucket που τις αντιστοιχεί. Οι συναρτήσεις που υλοποιούνται στην κλάση LSH είναι οι εξής.

- **nearest_neighbor**

Αυτή η συνάρτηση δέχεται σαν όρισμα ένα vector με το query και μας επιστρέφει τον ένα pair που περιέχει τον κοντινότερο γείτονα μαζί με την απόσταση που έχει από αυτό τον γείτονα. Η διαδικασία αυτή γίνεται υπολογίζοντας αρχικά το bucket που αντιστοιχεί στο query σε κάθε hashtable και στην συνέχεια πέρνουμε όλα τα στοιχεία που βρίσκονται σε αυτό το bucket στο vector image_indexes. Αφού αποθηκεύσουμε στο img_indexes προσορινά του κοντινούς γείτονες βρίσκουμε την Μανχάταν απόσταση μεταξύ αυτών και του query. Τέλος πέρνουμε την πιο κοντινή απόσταση από όλα και την επιστέφουμε

- **knn**

Η συνάρτηση αυτή λειτουργεί με παρόμοιο τρόπο με την nearest_neighbor με την μόνη διαφορά αντι να επιστρέψει ένα κοντινό γείτονα επιστρέφει του

K κοντινούς γείτονες. Οπότε εδώ δέχετε σαν όρισμα το query και το K που μας προσδιορίζει τον αριθμό των γειτόνων που θέλουμε να επιστρέψουμε με αποτέλεσμα να επιστρέφει ένα vector με K pairs που περιέχουν τον την απόσταση και τον N κοντινότερο γείτονα του query

- **range_search**

Η συνάρτηση range_search βρίσκει τους γείτονες του query απόσταση r. Δέχετε σαν όρισμα το query, την ακτίνα του κύκλου όπου θα γίνει το range search και μια σταθερά c, όπου αν δεν δόσουμε όρισμα πέρνει default τιμή 1. Στην συνέχεια όπως και οι προηγούμενες συνάρτησεις έτσι και η range search βρίσκει το bucket που αντιστοιχεί στο query σε κάθε hashtable και στην συνέχεια πέρνουμε όλα τα στοιχεία που βρίσκονται σε αυτό το bucket στο vector image_indexes. Τώρα για κάθε εικόνα ελέγχει βάση τις μετρικής Μανχάταν αν βρίσκετε εντός της ακτίνας r. Σε περίπτωση που βρίσκετε τότε προσθέτει την εικόνα στο results έτσι ώστε στο τέλος να τις επιστρέψει

- **exact_nearest_neighbor**

Αυτή η συνάρτηση έχει σκοπό να μας επιστρέψει τον ακριβές πιο κοντινους γείτονες για να ελέξουμε ότι τα αποτελέσματα των παραπάνω συναρτήσεων είναι σωστά. Η διαδικασία αυτή γίνεται με την μέθοδο του brute force, δηλαδή ελέγχουμε όλες τις εικόνες του dataset και επιστρέφουμε τις K εικόνες με την μικρότερη Μανχάταν απόσταση από το query. Εδώ η συνάρτηση αυτή πέρνει σαν όρισμα το query, το K μα επιστέφει ένα vector με pairs όπου το κάθε ζευγάρι αποτελείται από την απόσταση και στον εικόνα που είναι πιο κόντα στο query

3.5 HYPER CUBE

Για την υλοποίηση του Hyper Cube δημιουργήσαμε μια κλάση με το όνομα BinaryHyperCube η οποία έχει σαν σταθερές το d, M, probes, R καθώς και τρία vectors οποία είναι τα δεδομένα των εικόνων, οι τιμές των s και μια δομή για τον υπερκύβο. Στον **constructor** του υπερκύβου αρχικοποιούνε όλες η μεταβλητές και μπένουν τα

δεδομένα στο data vector. Στην συνέχεια δημιουργούνται με τυχαίο τρόπο τα s και μπένουν στο αντιστοιχο vector. Τέλος τα δεδομένα hasharοντε και μπένουν στο ανάλογο bucket της δομής του υπερκύβου. Η κλάση BinaryHyperCube υλοποιεί και τις παρακάτω συναρτήσεις.

- **f**

Η συνάρτηση F σύμφωνα με την θεωρία για την υλοποίηση του αλγοριθμου του υπερκύβου πρέπει να επιστρέφει 0 ή 1 με ομοιόμορφη κατανομή. Έτσι λοιπόν αποφασίσαμε η συνάρτηση f να δέχεται ένα integer, να πέρνει την διαδική του μορφή και να μετράει πόσα μηδινικά και πόσους άσσους έχει. Αν οι ασσοι είναι περισσότεροι από τα μηδινικά τότε επιστρέψει 1 αλλιώς επιστρέψει 0.

- **h**

Η συνάρτηση h είναι η hashfunction που χρησιμοποιά ο υπερκύβος αλλά είναι ίδια με την συνάρτηση $h(p)$ του Lsh που αναφέραμε πιο πάνω. Οπότε η υλοποίηση είναι η ίδια

- **get_number_from_bits**

- **hamming_distance**

- **knn**

- **range_search**

- **exact_nearest_neighbor**

- **get_bucket_imgs**

3.6 CLUSTERING

3.6.1 LLOYDS

3.6.2 LSH RANGE SEARCH