



ΕΛΛΗΝΙΚΗ ΔΗΜΟΚΡΑΤΙΑ

Εθνικόν και Καποδιστριακόν  
Πανεπιστήμιον Αθηνών

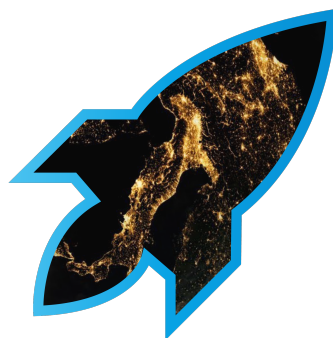
— ΙΔΡΥΘΕΝ ΤΟ 1837 —

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

<b>1</b>	<b>ΕΙΣΑΓΩΓΗ</b>	<b>3</b>
<b>2</b>	<b>ΜΕΤΑΓΛΩΤΤΙΣΗ-ΕΚΤΕΛΕΣΗ</b>	<b>4</b>
<b>3</b>	<b>ΥΛΟΠΟΙΗΣΗ</b>	<b>5</b>
3.1	ΕΙΣΟΔΟΣ ΔΕΔΟΜΕΝΩΝ . . . . .	5
3.2	ΜΕΤΡΙΚΕΣ . . . . .	6
3.3	LSH . . . . .	6
3.4	HYPER CUBE . . . . .	6
3.5	CLUSTERING . . . . .	6
3.5.1	Lloyds . . . . .	6
3.5.2	LSH Range Search . . . . .	6

---

## 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 METRIKES

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 απόσταση τους καθώς και ακόμα ένα όρισμα το οποίο είναι το όνομα της μετρική π.χ. L1 για την μετρική Μαχνάταν. Έδω μπορούν να υλοποιηθούν και άλλες μετρικές αλλά στην εργασία μας ζητήθηκες μόνο η μετρική Μανχάταν. Για την υλοποίηση της Μαχάταν μετρικής πήραμε το αθροισμα της απόλυτη τιμή των σημείων των δύο εικόνων

## 3.3 LSH

## 3.4 HYPER CUBE

## 3.5 CLUSTERING

### 3.5.1 LLOYDS

### 3.5.2 LSH RANGE SEARCH