

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Όνοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

## Εισαγωγή

- Η εργασία αναπτύχθηκε σε περιβάλλον **ubuntu 16.04.3**
- Χρησιμοποιήθηκε η **STL** της **C++ 11 (-std=c++11)** .
- Απο STL containers έχει χρησιμοποιηθεί μόνο:
  1. Vector
  2. Unordered\_set.
  3. Unordered\_map.
- Για τους χρόνους χρησιμοποιήθηκε η βιβλιοθήκη **chrono**.
- Υπάρχει η δυνατότητα μεταγλώττισης με Makefile. Με την εντολή make δημιουργείται 1 εκτελέσιμο:
  1. ./cluster
- Η εκτέλεση γίνεται όπως τις οδηγίες της άσκησης :
  1. **./cluster -i <input file> -c <configuration file> -o <output file> -d <metric>**
- Όλα τα ορίσματα πρέπει να δωθούν.
- Σε περίπτωση κάποιου λάθους με τα ορίσματα εμφανίζονται όλα τα λάθοι που εντοπίστηκαν και το πρόγραμμα τερματίζει.
- Στο αρχείο param.h υπάρχουν σταθερες και τύποι μεταβλητών που μπορούν να αλλάζουν ευκολα για τις ανάγκες της εργασίας .
- Χρησιμοποιήθηκε github με gitkraken ( [https://github.com/d1rafaelol/project\\_2.git](https://github.com/d1rafaelol/project_2.git)).

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Όνοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

## Αρχεία

Μερικές από τις συναρτήσεις είναι ίδιες σε αρκετά αρχεία. Αυτές θα περιγράφονται μόνο μια φορά.

### ■ check\_info.h

#### ◆ void print\_err(int )

Εκτυπώνει μήνυμα λάθους ανάλογα με τον αριθμό .

#### ◆ int check\_info\_cluster(int ,char\*\*,std::ifstream \*,std::ifstream \*,std::ofstream \*,std::string \*)

Ελέγχει τα ορίσματα ,ανοίγει τα αρχεία και δίνει τιμές σε μεταβλητές του προγράμματος .

### ■ Eucl\_Cos.h

#### ◆ Euclidean(int ,int );

Δίνει τυχαίες τιμές στα coordinates του διανύσματος vector\_V και στο t.

#### ◆ int Euclidean ::h\_func(Vector & )

Υπολογίζει το h(p) από την σελίδα 16 στο nnLSH.pdf.

#### ◆ Cosine(int)

Δίνει τυχαίες τιμές στα coordinates του διανύσματος vector\_R.

#### ◆ int Cosine ::h\_func(Vector & )

Υπολογίζει το h(p) από την σελίδα 22 στο nnLSH.pdf.

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Ονοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

### ■ General\_functions.h

#### ◆ int mod (long long int , unsigned long long int )

Δίνει το σωστό modulo ,όπως έχει αναφερθεί εδώ(  
[https://eclass.uoa.gr/modules/forum/viewtopic.php?course=DI352&topic=20690&forum=30954&fbclid=IwAR1U1ah0d6sLyAuB\\_a-Wnwplp7b4Wkv8Yal1Kf-lyKp7EAJ2QI11u49rlf0](https://eclass.uoa.gr/modules/forum/viewtopic.php?course=DI352&topic=20690&forum=30954&fbclid=IwAR1U1ah0d6sLyAuB_a-Wnwplp7b4Wkv8Yal1Kf-lyKp7EAJ2QI11u49rlf0) ).

#### ◆ void read\_info(std::ifstream \*,std::vector<Vector\*> \*)

Διαβάζει το αρχείο csv με τα διανύσματα και τα προσθέτει όλα σε ένα vector.

#### ◆ int check\_parameters(std::ifstream \*,int \*,int \*,int \*,int\*,int\* )

Τσεκάρει όλες τις παραμέτρους,και τους δίνει τιμές.(όλες έχουν και μια default τιμή σε περίπτωση που δεν δοθεί καποια, εκτός από τον αριθμό των clusters)

#### ◆ int check\_dimensions(std::vector<Vector\*> & )

Τσεκάρει αν όλα τα διανύσματα έχουν τον ίδιο αριθμό διαστάσεων.

#### ◆ int check\_cluster\_num(std::vector<Vector\*> & ,int)

Ελέγχει τον αριθμό των clusters να μην είναι πάνω από ένα οριο.Το οποίο μπορεί να αλλάξει από την μεταβλητή CLUSTERS\_DIV στο param.h.

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Όνοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

### ■ Hash.h

#### ◆ int Euclidean\_Hash::hashFunction(Vector \*)

Υπολογίζει το  $\phi(p)$  από την σελίδα 19 στο nnLSH.pdf.

#### ◆ int Euclidean\_Hash::g\_fun(Vector \*, Vector \*)

Ελέγχει αν δυο διανύσματα θα έπρεπε να είναι στο ίδιο bucket με την hashfunction  $g(p)$  στην σελίδα 19 στο nnLSH.pdf.

Υπολογίζει την  $h()$  για κάθε διάνυσμα και κοιτάει για ισότητες αν βρει τόσες όσες και οι  $k$  διαφορετικές  $h()$  τότε θα έμπαιναν στο ίδιο bucket οπότε επιστρέφει 1, διαφορετικά επιστρέφει 0

#### ◆ int Cosine\_Hash::hashFunction(Vector \*)

Υπολογίζει το  $\phi(p)$  από την σελίδα 23 στο nnLSH.pdf. Η  $g(x)$  είναι ένα int το οποίο είναι ίσο με το δυαδικό αριθμό  $\text{pch}(h_1(x)=1, h_2(x)=0, h_3(x)=1 \dots 101=5)$

#### ◆ int Cube\_Hash::hashFunction(Vector \*)

Βρίσκει την τιμή για κάθε  $h()$ . Αν αυτή η τιμή υπάρχει στο map τότε περνουμε την αντίστοιχη τιμή (0 1) από εκεί αν όχι την υπολογίζει και την προσθέτει στο map. Αυτές οι τιμές δημιουργούν έναν αριθμό (από 0 έως  $2^k-1$ ) ο οποίος είναι το bucket που πρέπει το διάνυσμα να μπει.

#### ◆ void insert Vector(Vector\*)

Βάζει ένα διάνυσμα στο hash table στο κατάλληλο bucket χρησιμοποιώντας την hashFunction.

#### ◆ std::vector<Vector\*> \* get\_similar\_Vectors(Vector\*)

Επιστρέφει το bucket που θα έμπαινε το διάνυσμα.

#### ◆ int get\_bucket\_num()

Επιστρέφει τον αριθμό των buckets.

#### ◆ std::vector<Vector\*> \* get\_bucket(int)

Επιστρέφει bucket ανάλογα με τον αριθμό.

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Ονοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

### ■ Neighbour Namespaces.h Distance.h (Αλλαγή απο 1o project).

#### Λειτουργεί με κλάση αντί για namespace

- ◆ void Range\_N(std::vector<std::string> &, std::vector<Vector\*> &, Vector \*, double)

Βάζει σε vector όλα τα διανύσματα που έχουν απόσταση λιγότερο από R το οποίο δίνεται στο αρχείο με τα queries.

- ◆ void Range\_bet\_N(std::vector<Vector\*> &, std::vector<Vector\*> &, Vector \*, long double, long double)

Βάζει σε vector όλα τα διανύσματα που έχουν απόσταση ανάμεσα στο σε δύο ακτίνες

- ◆ void Nearest\_N(std::vector<std::pair<std::string, double>> &, std::vector<Vector\*> &, Vector \*)

Βάζει σε vector όλα τα διάνυσμα με την μικρότερη απόσταση από το διάνυσμα του query

Η απόσταση υπολογίζεται με διαφορετικό τρόπο ανάλογα με το namespace.

#### double cosi(Vector &, Vector &)

σελίδα 21 στο nnLSH.pdf. Για τον υπολογισμό της απόστασης cosine  $1 - \text{cosi}$

#### double dist(Vector & x, Vector & y)

σελίδα 14 στο nnLSH.pdf. Για  $k=2$

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Ονοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

### ■ Vector.h

#### ◆ Vector(std::string )

Δίνει τιμή στο id του διανύσματος.

#### ◆ void add\_coordinate(coordinate )

Προσθέτει ένα καινούργιο coordinate στο vector του διανύσματος.

#### ◆ std::string get\_identity(void)

Επιστρέφει το id του διανύσματος.

#### ◆ std::vector<coordinate> \*get\_coordinates(void)

Επιστρέφει τις συντεταγμένες του διανύσματος.

### ■ Cluster.h

#### ◆ void add\_Vector(Vector \*)

Βάζει ένα Vector στο cluster.

#### ◆ int update\_Centroid(Vector \*)

Αλλάζει το centroid με τον αλγόριθμο k-means.

#### ◆ std::vector<Vector\*> \*get\_cluster\_vectors()

Επιστρέφει ένα vector με τα διανύσματα που είναι μέσα στο cluster εκτός του centroid.

### ■ Cluster\_Group.h

#### ◆ void k\_unique\_rand\_init(std::vector<Vector\*> & )

Παίρνει τυχαία k clusters από τα διαθέσιμα vectors.

#### ◆ void k\_means\_plus\_plus\_init(std::vector<Vector\*> & )

Παίρνει k clusters από τα διαθέσιμα vectors σύμφωνα με το k-means++.

# Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

## Project 2

Ονοματεπώνυμο: Χατζηδάκης Ραφαήλ

Αριθμός Μητρώου: 1115201400248

- ◆ **void k\_means\_average\_init(std::vector<Vector\*> &)**

Παιρνει τυχαία ένα cluster και μετά υπολογίζει την μέση τιμή των αποστάσεων όλων των μέχρι τώρα επιλεγμένων clusters και παίρνει με πιθανότητα όπως και στο k-means\_plus\_plus αυτο με την μεγαλύτερη μέση τιμή για επόμενο centroid.

- ◆ **void Lloyd\_assignment(std::vector<Vector\*> &)**

Βάζει όλα τα διανύσματα σε κάποιο cluster σύμφωνα με τον αλγόριθμο Lloyd's .

- ◆ **void Range\_search\_assignment\_LSH(std::vector<Vector\*> &, int ,int );**

- ◆ **void Range\_search\_assignment\_Hypercube(std::vector<Vector\*> &,int,int,int);**

Range search κάνουν το ίδιο αλλά με διαφορετικό hash table και με την διαφορά ότι το hypercube κοιτάει και τα probes. Λειτουργουν σύμφωνα με τον αλγόριθμο range search. Στην αρχή ψάχνουμε όλα τα vectors που πρέπει να ελεγχθούν απο range search από κάθε hash table και κοιτάμε για ίδια διανύσματα στην αρχή ,μετά ο αλγόριθμος συνεχίζεται κανονικά.

- ◆ **int k\_means\_update()**

Δημιουργεί k καινούργια clusters που δεν ανήκουν στο αρχείο εισόδου σύμφωνα με το αλγόριθμο k-means.

- ◆ **int PAM\_improvement\_update()**

Αλλάζει τα clusters με τα medoids .

Το αρχείο cluster.conf παιρνει 2 παραπάνω παραμετρους το 1) Number of probes και 2) number\_of\_vectors. Τα οποία παιρνουν default τιμες απο το param.h αν δεν δωθουν στο αρχείο.