

Κ23γ: Ανάπτυξη Λογισμικού για Αλγοριθμικά Προβλήματα

Χειμερινό εξάμηνο 2014-15

2η Προγραμματιστική Εργασία

Υλοποίηση του αλγορίθμου ομαδοποίησης K-medoids στη γλώσσα C

Α) Στοιχεία φοιτητών που ανέπτυξαν το πρόγραμμα.

<u>Επώνυμο:</u>	Μαντάς	Μεγκραμπιάν
<u>Όνομα:</u>	Ιωάννης	Αρσάκ
<u>A.M.:</u>	1115200700298	1115201100049
<u>email:</u>	sdi0700298@di.uoa.gr	sdi1100049@di.uoa.gr

Β) Τίτλος και περιγραφή του προγράμματος

Ο τίτλος του προγράμματος είναι SDproj2.

Λειτουργικό Σύστημα Υλοποίησης: Ubuntu 14.04

IDE Υλοποίησης: Geany v.1.23.1

Χρόνος Υλοποίησης: 3 εβδομάδες

Μέθοδοι Αρχικοποίησης (Initialization)

- **Concentrate:** Ο αλγόριθμος αυτός υλοποιήθηκε όπως αναφέρεται στο “*A simple and fast algorithm for K-medoids clustering*” [Park-Jun '09], του οποίου το implementation αναφέρεται και στις διαφάνειες του μαθήματος. Το πρόγραμμα μας υλοποιεί πλήρως τον αλγόριθμο αυτό.
- **K-medoids++:** Ο αλγόριθμος αυτός υλοποιήθηκε όπως αναφέρεται στο “*K-means++: The Advantages of Careful Seeding*” [Arthur - Vassilvitskii '07], του οποίου το implementation αναφέρεται και στις διαφάνειες του μαθήματος. Το πρόγραμμα μας υλοποιεί πλήρως τον αλγόριθμο αυτό.

Μέθοδοι Ανάθεσης (Assignment)

- **PAM assignment (simplest approach):** Ο αλγόριθμος αυτός για το assignment, όπως αναφέρει και το όνομα του, αποτελεί την απλούστερη μορφή για ανάθεση των αντικειμένων του συνόλου στα κέντρα των clusters. Δεν χρησιμοποιεί κανενός είδους βελτίωση απλώς ελέγχει όλα τα σημεία με τα κέντρα εξαντλητικά, όπως μας δίνεται και στις διαφάνειες.
- **Assignment by LSH/DBH (Reverse Approach):** (Edit: Δεν προλάβουμε τελικά να υλοποιήσουμε να κάνει compile και να είναι πλήρως συμβατός με το πρόγραμμα το LSH assignment.) Μπορούμε εφόσον μας δοθεί η δυνατότητα να παραδώσουμε εντός ημερών, καθώς βρισκόμασταν πολύ

κοντά στην ολοκλήρωση).

Μέθοδοι Ενημέρωσης (Update)

- **“A la Lloyd's”(improved PAM):** Ο αλγόριθμος αυτός για το update, όπως αναφέρεται στο *“A simple and fast algorithm for K-medoids clustering” [Park-Jun '09]*, του οποίου το implementation αναφέρεται και στις διαφάνειες του μαθήματος αλλά και συζητήθηκε στο φροντιστήριο Το πρόγραμμα μας υλοποιεί πλήρως τον αλγόριθμο αυτό.
- **CLARANS:** Ο αλγόριθμος αυτός υλοποιήθηκε όπως ακριβώς αναφέρεται στο *“CLARANS: A Method for Clustering Objects for Spatial Data Mining” [Ng - Han '94]*. Δηλαδή χρησιμοποιεί το cost differentials και κάνει τα ανάλογα swaps μέσω του αλγορίθμου PAM. Λόγω αδυναμίας για υπολογισμού του 2nd best centroid μέσω του LSH assignment, δεν είμαστε σε θέση να συνδυάσουμε αυτές τις δύο μεθόδους.

Μέθοδοι Αποτίμησης (Evaluation)

- **Silhouettes(simplest approach):** Ο αλγόριθμος αυτός για το evaluation του clustering, υλοποιήθηκε όπως αναφέρεται στο *“Silhouettes: a graphical aid to the interpretation and validation of cluster analysis” [Rousseeuw '87]*. Ο αλγόριθμος
- Σε κάθε μέθοδο επίσης μετράται ο χρόνος, όπως ζητείται, για να κάνουμε τους αντίστοιχους ελέγχους.
- Για το LSH πήραμε τα “χρήσιμα” κομμάτια κώδικα από την προηγούμενη άσκηση και τα χρησιμοποιήσαμε για το LSH assignment.

Γ)Κατάλογος των αρχείων κώδικα /κεφαλίδων και περιγραφή τους.

Τα αρχεία έχουν χωριστεί σε λογικές ομάδες για λόγους επεκτασιμότητας, απόκρυψης δεδομένων, καλύτερης διαχείρισης των αρχείων αλλά για λόγους καλύτερης αντίληψης των προγραμμάτων.

Πηγαία Αρχεία

./main.c	Τρέχει το πρόγραμμα και καλεί τις ανάλογες συναρτήσεις
./aEucMain.c	Η κύρια συνάρτηση clustering για Euclidean Spaces
./aHamMain.c	Η κύρια συνάρτηση clustering για Hamming Distance.
./aMatMain.c	Η κύρια συνάρτηση clustering για Distance Matrixes.
./cluster.c	Συναρτήσεις για τα cluster.
./generalfuncs.c	Περιέχει διάφορες γενικές συναρτήσεις
./algorithms.c	Περιέχει τις υλοποιήσεις των ζητούμενων αλγόριθμων.

<code>./Euc/eucfuncs.c</code>	Περιέχονται συναρτήσεις σχετικές με <code>aEucMain.c</code>
<code>./Ham/hamfuncs.c</code>	Περιέχονται συναρτήσεις σχετικές με <code>aHamMain.c</code>
<code>./Mat/matfuncs.c</code>	Περιέχονται συναρτήσεις σχετικές με <code>aMatMain.c</code>



Αρχία Κεφαλίδων

<code>./headers.h</code>	Περιέχει όλες τα υπόλοιπα αρχία κεφαλίδων
<code>./mainfuncs.h</code>	Περιέχει τα ορίσματα των 3 κύριων συναρτήσεων.
<code>./cluster.h</code>	Τα ορίσματα του αντίστοιχου <code>random.c</code>
<code>./defines.h</code>	Περιέχει τα <code>defines</code> του προγράμματος.
<code>./generalfuns.h</code>	Τα ορίσματα του αντίστοιχου <code>generalfuns.c</code>
<code>./algorithms.h</code>	Τα ορίσματα του αντίστοιχου <code>algorithms.c</code>
<code>./Euc/eucfuncs.h</code>	Τα ορίσματα του αντίστοιχου <code>aEucMain.c</code>
<code>./Ham/hamfuncs.h</code>	Τα ορίσματα του αντίστοιχου <code>aHamMain.c</code>
<code>./Mat/matfuncs.h</code>	Τα ορίσματα του αντίστοιχου <code>aMatMain.c</code>

Λοιπά

<code>./makefile</code>	Το αρχείο <code>makefile</code> .
<code>./cluster.conf</code>	Το <code>configurator file</code> .
<code>./LSH</code>	Directory που περιέχει όλα τα αρχία για το LSH/DBH.

Δ)Οδηγίες μεταγλώττισης του προγράμματος.

Εντολές του `makefile`

make	Δημιουργεί όλα τα αντικειμενικά αρχία και το εκτελέσιμο.
make clean	Διαγράφει όλα τα αντικειμενικά αρχία και το εκτελέσιμο.
make count	Μετράει πλήθος γραμμών, λέξεων και χαρακτήρων.

Ε)Οδηγίες χρήσης του προγράμματος

Το πρόγραμμα εκτελείται γράφοντας

`./SDproj1 -d <input file> -c <configuration file> -o <output file> -complete`

Όλα τα ορίσματα είναι προαιρετικά. Σε περίπτωση που δεν δοθούν τα L και k έχουν default τιμές ενώ τα διάφορα files ζητούνται στην συνέχεια.

Ο χρήστης θα δώσει ένα dataset, σαν τα ενδεικτικά που μας δόθηκαν και το πρόγραμμα θα του επιστρέψει σε ένα αρχείο για όλους τους συνδυασμούς των παραπάνω αλγόριθμων που αναφέρθηκαν ό,τι ακριβώς ζητείται στην εκφώνηση της άσκησης.