

Nearest Items and Clustering

Πίνακας Περιεχομένων

[Πίνακας Περιεχομένων](#)

[Στοιχεία](#)

[Οδηγίες μεταγλώττισης και εκτέλεσης του προγράμματος](#)

[Περιγραφή](#)

[Κατάλογος των αρχείων](#)

[LSH Αρχεία και Περιγραφή](#)

[Hypercube Αρχεία και Περιγραφή](#)

[Cluster Αρχεία και Περιγραφή](#)

[k-Means++ Αρχικοποίηση](#)

[Συνθήκες Σύγκλισης](#)

[Βοηθητικές Συναρτήσεις](#)

Στοιχεία

ΝΤΑΡΟΥΙΣ ΝΙΖΑΡ 1115201800286
ΤΣΙΛΙΦΩΝΗΣ ΑΡΗΣ 1115201700170

Οδηγίες μεταγλώττισης και εκτέλεσης του προγράμματος

Για την μεταγλώττιση έχει φτιαχτεί ένα makefile το οποίο φτιάχνει τα απαραίτητα εκτελέσιμα για όλα τα προγράμματα.

Η εκτέλεση των προγραμμάτων γίνεται ακριβώς όπως ζητήθηκε στην εκφώνηση. Οι παράμετροι γράφονται με οποιαδήποτε σειρά.

Περιγραφή

Η εργασία υλοποιεί τους αλγορίθμους αναζήτησης κοντινότερου γείτονα, LSH και Hypercube χρησιμοποιώντας την μετρική L2 (euclidean distance). Καθώς και τους αλγορίθμους συσταδοποίησης διανυσμάτων στον χώρο, Lloyds και reverse Range search χρησιμοποιώντας τον LSH και Hypercube, με την ίδια μετρική. Τέλος, η αρχικοποίηση των διανυσμάτων που θα χρησιμοποιηθούν από τους αλγόριθμους συσταδοποίησης πραγματοποιείται με την τεχνική k-Means++.

Όλα τα ζητούμε της άσκησης έχουν υλοποιηθεί.

Το πρόγραμμα δεν αποθηκεύει διπλότυπα δεδομένα, στα hash tables αποθηκεύονται τα Item_id.

Κατάλογος των αρχείων

LSH Αρχεία και Περιγραφή

Τα LSH αρχεία θα τα βρείτε στον κατάλογο

“Algorithmic-Problems-Project1\algo\LSH_Folder” και το πρόγραμμα (main) που εκτελεί τους μεθόδους και τον αλγόριθμο βρίσκεται στον κατάλογο

“Algorithmic-Problems-Project1\algo”.

Στο αρχείο “LSH/LSH.hpp” θα βρείτε τις δομές, τις μεταβλητές και τις μεθόδους που χρησιμοποιούνται στον αλγόριθμο, καθώς και σχόλια εξηγώντας τι κάνει η κάθε μέθοδο, δομή και σε τι χρησιμεύει η κάθε μεταβλητή.

Η συνάρτηση κατακερματισμού, και η τιμή της, υπολογίζεται στην συνάρτηση “Specific Hash Value” και έχει υλοποιηθεί ακριβώς έτσι όπως έλεγε η θεωρεία στις διαφάνειες. Το w ορίζεται με τον υπολογισμό της μέσης ευκλείδειας απόστασης μεταξύ των 5% των δεδομένων εισόδου, εάν η είσοδος είναι μικρή και επομένως το 5% είναι μικρότερο από το μηδέν, παίρνουμε το 50% των δεδομένων. Η τιμή που

παίρνει το w πρέπει να εξαρτάται από τα δεδομένα. Καταλήξαμε στο 5% των δεδομένων πειραματίζοντας με διάφορα ποσοστά και παρατηρώντας τα αποτελέσματα.

Όσον αφορά το πλήθος των buckets παρατηρήσαμε πως το $n/4$ είναι αρκετά καλή τιμή καθώς και είναι πιο σκορπισμένα τα δεδομένα με αποτέλεσμα ο αλγόριθμος να είναι πιο γρήγορος.

Οι αλγόριθμοι search by range και approximate nearest neighbors έχουν υλοποιηθεί στο αρχείο "LSH_Folder/LSH.cpp" με πλήρη περιγραφή για την λειτουργία τους, πάνω από τις συναρτήσεις "Search_by_range" και "Nearest_N_search".

Hypercube Αρχεία και Περιγραφή

Τα Hypercube αρχεία θα τα βρείτε στον κατάλογο

"Algorithmic-Problems-Project1\algo\Hypercube" και το πρόγραμμα (main) που εκτελεί τους μεθόδους και τον αλγόριθμο βρίσκεται στον κατάλογο "Algorithmic-Problems-Project1\algo".

Στο αρχείο "Hypercube/Hypercube.hpp" θα βρείτε τις δομές, τις μεταβλητές και τις μεθόδους που χρησιμοποιούνται στον αλγόριθμο, καθώς και σχόλια εξηγώντας τι κάνει η κάθε μέθοδο, δομή και σε τι χρησιμεύει η κάθε μεταβλητή.

Αρχικά έχω δημιουργήσει ένα πίνακα όπου για κάθε θέση έχω ζευγάρια (h_i, f_i) . Έτσι μπορώ να αποθηκεύω για κάθε τιμή του $H[i]$ την αντίστοιχη τιμή f_i (το $H[1]$ έχει πολλά $\{h_1, f_1\}$). Δηλαδή για $h_1(\text{point}_1) = 3$ και $f_1(h_1) = 0$ σώζω το ζευγάρι $\{h_1, f_1\}$ ώστε άμα ξανατύχει κάποιο point_x να έχει $h_1(\text{point}_x) = 3$ να μην χρειαστεί να ξαναυπολογίσω το $f_1(h_1)$. Έτσι δημιουργείται μια δομή hashtable όπου για κάθε $H[i]$ βρίσκω ή υπολογίζω την αντίστοιχη τιμή f_i . Με αυτόν τον τρόπο δημιουργείται για κάθε σημείο p η εξής συμβολοσειρά : $[f_1(h_1(p)), f_2(h_2(p)) \dots, f_d(h_d(p))]$, $d' = k$ στην υλοποίηση. Στην συνέχεια, αυτή η συμβολοσειρά που προκύπτει (πχ. $k = 4 \rightarrow 0110$) είναι η δυαδική αναπαράσταση ενός δεκαδικού αριθμού, που αντιστοιχεί στην θέση του bucket στο hypercube hashtable. Οπότε $0110 = 7$, συνεπώς το σημείο μπαίνει στον bucket 7. Όσον αφορά την αναζήτηση των queries, πρώτα βρίσκω το bucket που βάζει ο αλγόριθμος hypercube το query. Έπειτα βρίσκω με hamming distance τους γειτονικούς bucket σε αυτόν τον bucket. Δηλαδή για κάθε bucket βρίσκω το hamming distance του bucket με το bucket που έπεσε το query. Όπως γνωρίζουμε το hamming distance βρίσκει πόσους διαφορετικούς άσους και μηδενικά έχουν δύο ακέραιοι. (σαν να σου λέει πόσα διαφορετικά f_i έχουν). Αφού ταξινομήσω τα γειτονικά buckets ανάλογα με την απόσταση του query bucket, ψάχνω τόσους buckets όσους μου λέει η παράμετρος probes. Για κάθε bucket, ψάχνω όλα τα σημεία του και αναζητώ πόσα είναι μέσα στο radius που έχω δώσει σαν παράμετρο στην συνάρτηση. Όσα είναι εντός, επιστρέφω το γειτονικό σημείο σε μια λίστα και την απόσταση του στην δεύτερη λίστα. Αντίστοιχα στο N-nearest neighbours απλά καθώς ψάχνω τα probes (buckets) έχω έναν vector όπου τον γεμίζω με γείτονες και τον ταξινομώ ώστε ο μακρινότερος να είναι στο τέλος. Αν ο vector με τα nearest

neighbors γεμίσει και έχω βρει σημείο για υποψήφιο κοντινό , κοιτάω τον τελευταίο γείτονα .Σε περίπτωση που είναι πιο κοντά, πετάω τον παλιό και βάζω στην θέση του τον καινούριο.Έπειτα ταξινομώ και πάλι το ίδιο.(Κάτι σαν min heap).

Cluster Αρχεία και Περιγραφή

Τα Cluster αρχεία θα τα βρείτε στον κατάλογο

“Algorithmic-Problems-Project1\algo\kMeans++” και το πρόγραμμα (Cluster_main.cpp) που εκτελεί τις μεθόδους και τον αλγόριθμο βρίσκεται στον κατάλογο “Algorithmic-Problems-Project1\algo”.

Στο αρχείο “kMeans++/Cluster.hpp” θα βρείτε τις δομές, τις μεταβλητές και τις μεθόδους που χρησιμοποιούνται στον αλγόριθμο, καθώς και σχόλια εξηγώντας τι κάνει η κάθε μέθοδο, δομή και σε τι χρησιμεύει η κάθε μεταβλητή.

Στο αρχείο “kMeans++/Cluster.cpp” θα βρείτε την υλοποίηση των μεθόδων kMeans++ (συνάρτηση kMeanspp_Initialization), Lloyd (συνάρτηση Lloyd_method), Silhouette (συνάρτηση Silhouette) και Reverse search (συνάρτηση reverse_assignment), καθώς και την πλήρη περιγραφή τους πάνω από κάθε συνάρτηση.

Όλα υλοποιήθηκαν όπως στις διαφάνειες.

k-Means++ Αρχικοποίηση

Η μέθοδος έχει υλοποιηθεί ως εξής:

- Το πρώτο κεντροειδές επιλέγεται ομοιόμορφα
- Τα υπόλοιπα κεντροειδή επιλέγονται με βάση την πιθανότητα (Πιθανότητα = μικρότερη απόσταση/άθροισμα μικρότερων αποστάσεων)
- μικρότερη απόσταση = η απόσταση του σημείου σε σχέση με το πλησιέστερο μικρότερων κεντροειδή
- Οι υπολογισμοί γίνονται με την ευκλείδεια απόσταση
- Το σημείο με την μεγαλύτερη πιθανότητα (δηλαδή το πιο μακρινό από όλα τα είδη υπάρχοντα κεντροειδή) θα είναι το επόμενο κεντροειδή.

Συνθήκες Σύγκλισης

Ο αλγόριθμος Lloyd για να τερματίσει συγκρίνει όλα τα νέα κεντροειδή με τα προηγούμενα. Συγκεκριμένα συγκρίνει τις συντεταγμένες τους και αν έχουν αλλάξει λιγότερο από 1% (όλων των κεντροειδών), τότε τερματίζει.

Ο αλγόριθμος reverse τερματίζει όταν η απόσταση μεταξύ του νέου και του προηγούμενου κέντρου είναι μικρότερη από 0.1 τουλάχιστον στα μισά κεντροειδή.

Βοηθητικές Συναρτήσεις

Στο αρχείο “Algorithmic-Problems-Project1\algo\functions” υπάρχουν βοηθητικές συναρτήσεις που επαναχρησιμοποιούνται καθόλη την διάρκεια εκτέλεσης των προγραμμάτων.