

ΑΝΑΠΤΥΞΗ ΛΟΓΙΣΜΙΚΟΥ ΓΙΑ ΠΛΗΡΟΦΟΡΙΑΚΑ ΣΥΣΤΗΜΑΤΑ

ΤΕΛΙΚΗ ΑΝΑΦΟΡΑ 2023-2024

ΠΙΤΣΑΣ ΛΟΥΚΑΣ - 1115201900156

ΣΤΕΡΓΙΟΣ ΤΖΙΡΑΚΗΣ - 1115201900188

Εισαγωγή

Σε αυτή την αναφορά, παρουσιάζουμε το τελικό παραδοτέο στα πλαίσια του μαθήματος *Ανάπτυξη Λογισμικού για Πληροφοριακά Συστήματα* για το ακαδημαϊκό έτος 2023-2024. Το πρόβλημα προς επίλυση είναι η κατασκευή KNN γράφων για την εύρεση k-κοντινότερων γειτόνων πάνω σε πολυδιάστατα σύνολα δεδομένων. Η εφαρμογή που αναπτύξαμε σε C++ ακολουθεί την λογική του αλγορίθμου *NN-Descent* όπως ήταν ζητούμενο. Υλοποιούνται διάφορες βελτιστοποιήσεις που περιγράφονται στην υλοποίηση του [PyNNDescent](#) καθώς και από τους [Wei Dong, Charikar Moses, Kai Li, 2011](#).

Στην αναφορά μας, εξηγούμε τις σχεδιαστικές επιλογές μας για την ανάπτυξη του εν λόγω αλγορίθμου καθώς και των λοιπών βελτιστοποιήσεων. Παρουσιάζουμε πειραματικά δεδομένα και τα αποτελέσματα της εφαρμογής πάνω σε αυτά.

Δομές δεδομένων

Για τη κατασκευή του γράφου, υλοποιούμε τη παρακάτω δομή για να αναπαραστήσουμε ένα κόμβο μέσα στο σύνολο δεδομένων:

```
//A class defining the KNN (k-nearest neighbors) problem.
class KNN {
    //The data points become nodes inside our graph.
    //We define neighbors and reverse neighbors separately,
    //thus we have both a directed and an undirected graph at will.
    struct node {
        float norm{ 0 }; //Euclidean norm
        float *cord{ nullptr }; //Coordinate array

        omp_lock_t lock; //Lock for parallelism
        minAVL<float, unsigned int> Cedge; //Candidates found during solve
        N_AVL<float, unsigned int> edge; //Indices to neighbors, sorted by distance - only the k best are kept
        RN_AVL<unsigned int> Redge; //Indices to reverse neighbors
    };
};
```

Ο γράφος είναι απλά ένας πίνακας με n (graph size) από τις παραπάνω δομές. Οι κόμβοι μας εσωτερικά αποτελούνται σχεδόν εξ ολοκλήρου απο AVL δέντρα: οι γείτονες, οι πιθανοί γείτονες, καθώς και οι αντίστροφοι γείτονες, αναπαριστώνται με κάποια μορφή AVL δέντρων (των οποίων οι λεπτομέρειες υλοποίησης φαίνονται στο αρχείο AVL.h). Πρόκειται για templated δομές που δημιουργήσαμε απο την αρχή για υλοποίηση της εφαρμογής μας.

Η επιλογή AVL δέντρων βασίζεται στους εξής παράγοντες:

- **Οργάνωση Δεδομένων:** Τα AVL δέντρα που χρησιμοποιούνται για τους γείτονες ταξινομούνται αυτόματα. Με άλλα λόγια, οι γείτονες θα είναι πάντα ταξινομημένοι με βάση τη καλύτερη απόσταση, και δεν χρειάζεται παραπάνω ενέργεια για ταξινόμηση.

Οι γείτονες και οι αντίστροφοι γείτονες είναι χωρισμένοι σε νέους και παλιούς για την διαδικασία δειγματοληψίας. Η διάκριση γίνεται με τον διαχωρισμό σε δύο δέντρα, ένα για τους καινούργιους και ένα για τους παλιούς γείτονες. Αν και αυτό σημαίνει πως χρειάζονται μεταφορές μεταξύ των δύο δέντρων, θεωρούμε πως είναι μια καλή επιλογή, καθώς εξαιρείται εντελώς η ανάγκη αποθήκευσης επιπλέον πληροφορίας (bool flags), καθώς και η ανάγκη αναζήτησης των νέων γειτόνων κατά τη διαδικασία δειγματοληψίας.

Για χαμηλότερα ποσοστά δειγματοληψίας, η αποδοτικότητα της προσέγγισής μας συνεχώς αυξάνεται, καθώς μειώνονται οι μεταφορές, ενώ αυξάνεται ο λόγος: (χώρος που απαιτείται για αποθήκευση bool flags) / (πλήθος δειγμάτων).

- **Ταχύτητα:** Η εισαγωγή και εξαγωγή ενός στοιχείου στο δέντρο έχει χειρότερη απόδοση $O(\log n)$, ενώ το δέντρο παραμένει ισορροπημένο. Αυτό σημαίνει πως οποιαδήποτε αναζήτηση επίσης έχει χειρότερη απόδοση $O(\log n)$. Η αναζήτηση είναι απαραίτητη για την διαγραφή ενός αντίστροφου γείτονα.

Κατα το βήμα ενημέρωσης του KNN γράφου, ο χειρότερος γείτονας (edge) συγκρίνεται με τον καλύτερο πιθανό γείτονα (Cedge). Η εύρεση μέγιστου/ελάχιστου στα AVL δέντρα είναι εξαιρετικά γρήγορη, ενώ η όλη διαδικασία σύγκρισης σταματάει μόλις βρεθεί ο πρώτος πιθανός γείτονας που δεν προσφέρει βελτίωση, αφού και οι δύο δομές γειτόνων είναι ταξινομημένες.

- **Εξοικονόμηση χώρου:** Τα AVL δέντρα που χρησιμοποιούνται για τους γείτονες καθώς και για τους πιθανούς γείτονες (Cedges) επιτρέπουν την ύπαρξη μέχρι k μοναδικών αντικειμένων μέσα στη δομή. Δεν υπάρχουν ποτέ διπλότυπα, ενώ η προσπάθεια εισαγωγής νέου γείτονα ενώ τα δέντρα είναι γεμάτα θα επιτύχουν μόνο αν ο νέος γείτονας έρχεται με καλύτερη απόσταση από τον χειρότερο γείτονα.

Συνολικά η κλάση μας καταλαμβάνει όσο χώρο καταλαμβάνει ο γράφος, σύν μερικές επιπλέον πληροφορίες μικρού κόστους:

```
//The graph itself
vector<struct node> graph;

//Additional parameters
float *worst{ nullptr }; //Used on true solve
char dataset[100] = { '\0' };
unsigned int dim{ 0 }, threads{ 4 };
float (KNN::*distf)(unsigned int, unsigned int){ nullptr }; //Distance function
```

Στάδια εκτέλεσης

Για να είναι σαφής η δομή της εφαρμογής, αλλά και της κατανομής χρονισμού του προγράμματός μας, εξηγούμε συνοπτικά τα στάδια που ακολουθεί μια τυπική εκτέλεση του **KNN**:

- **Στάδιο 1 (initialization):** Κατα την κλήση της συνάρτησης `KNN::solve()` (αφού δηλαδή έχει διαβαστεί ο γράφος από το αρχείο `dataset`), ξεκινάει ο πρώτος χρονισμός του προγράμματός μας. Αυτό το στάδιο περιλαμβάνει την αρχικοποίηση των γειτόνων μέσω δέντρων τυχαίων προβολών (RPTs), και σε προηγούμενες εκδόσεις της εργασίας, μέσω k -τυχαίων γειτόνων.
- **Στάδιο 2 (solve):** Σε αυτό το στάδιο εμπίπτει η επαναληπτική διαδικασία βελτίωσης του KNN γράφου. Διαιρείται σε δύο φάσεις: την δημιουργία υποψηφίων γειτόνων για κάθε κόμβο (Region A), και έπειτα την ενημέρωση του γράφου με τους υποψηφίους (Region B). Στο πρόγραμμά μας η πρώτη φάση καταλαμβάνει το περισσότερο χρόνο, αφού σε αυτό γίνονται και οι υπολογισμοί αποστάσεων.
- **Στάδιο 3 (determine accuracy):** Με τη συνάρτηση `accuracy()`, εξετάζουμε την ακρίβεια του KNN γράφου με τον αυθεντικό (100% ευστοχία) γράφο. Εάν υπάρχει λυμένος ο γράφος για το αντίστοιχο k , η διαδικασία είναι γρήγορη καθώς οι πραγματικοί γείτονες εξάγονται από το αρχείο. Αν δεν υπάρχει τέτοιο αρχείο, οι γείτονες υπολογίζονται από την αρχή. Αυτή η διαδικασία είναι αρκετά χρονοβόρα.

Τα στάδια που ακολουθεί το εκτελέσιμο **GraphSolve** είναι παρόμοια.

```
File: 00010000-4.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 0 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 1905ms
Solving graph... graph solved at 12 iterations | Time elapsed: 13376ms
Calculating accuracy... relative accuracy: 99.6984% | Absolute accuracy: 99.4387% | Time elapsed: 672ms
```

*Τυπική εκτέλεση του **KNN**, 10000 nodes*

```
linux18:/home/users/sdi1900156/Documents/Project2023> ./GraphSolve 00005000-1.bin -k 50
Solving graph... graph solved | Time elapsed: 15426ms
Exporting graph... done | Time elapsed: 185ms
linux18:/home/users/sdi1900156/Documents/Project2023>
```

*Τυπική εκτέλεση του **GraphSolve**, 5000 nodes, $k=50$*

Τυχαία VS γρήγορη δειγματοληψία

Η διαγραφή εσωτερικών κόμβων του AVL δέντρου (δηλαδή εκείνων που δεν είναι κοντά στα φύλλα) είναι αρκετά πιο αργή από την διαγραφή φύλλων του δέντρου. Υπο κανονικές συνθήκες, η δειγματοληψία κατά το στάδιο εύρεσης πιθανών γειτόνων είναι ημι-τυχαία (χωρίς χρήση RNG), η οποία μας προσφέρει ιδανική ευστοχία στον γράφο (>99% για αντίστοιχο k).

Ωστόσο, υλοποιούμε και έναν πιο γρήγορο τρόπο δειγματοληψίας, ο οποίος επιφέρει κατά μέσο όρο 66% γρηγορότερη λύση του γράφου με περίπου 1% απώλεια στην ευστοχία. Η γρήγορη δειγματοληψία γίνεται με εξαγωγή των min κόμβων.

Παραδείγματα κανονικής VS γρήγορης δειγματοληψίας:

```
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00010000-4.bin
File: 00010000-4.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 0 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 1931ms
Solving graph... graph solved at 11 iterations | Time elapsed: 13268ms
Calculating accuracy... relative accuracy: 99.695% | Absolute accuracy: 99.4324% | Time elapsed: 656ms
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00010000-4.bin -f
File: 00010000-4.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 1918ms
Solving graph... graph solved at 9 iterations | Time elapsed: 5015ms
Calculating accuracy... relative accuracy: 98.6095% | Absolute accuracy: 97.5361% | Time elapsed: 764ms
```

VS fast sampling, 10000 nodes

```
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00002000-3.bin
File: 00002000-3.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 0 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 174ms
Solving graph... graph solved at 8 iterations | Time elapsed: 2239ms
Calculating accuracy... relative accuracy: 99.8853% | Absolute accuracy: 99.7609% | Time elapsed: 2454ms
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00002000-3.bin -f
File: 00002000-3.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 175ms
Solving graph... graph solved at 7 iterations | Time elapsed: 813ms
Calculating accuracy... relative accuracy: 99.4884% | Absolute accuracy: 98.7689% | Time elapsed: 2452ms
linux18:/home/users/sdi1900156/Documents/Project2023>
```

VS fast sampling, 20000 nodes

```
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00005000-1.bin
File: 00005000-1.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 0 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 657ms
Solving graph... graph solved at 10 iterations | Time elapsed: 6169ms
Calculating accuracy... relative accuracy: 99.781% | Absolute accuracy: 99.5735% | Time elapsed: 417ms
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00005000-1.bin -f
File: 00005000-1.bin | k: 20 | Delta: 0.001 | Sampling rate: 80% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 4 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 645ms
Solving graph... graph solved at 8 iterations | Time elapsed: 2298ms
Calculating accuracy... relative accuracy: 99.8451% | Absolute accuracy: 98.1476% | Time elapsed: 485ms
```

VS fast sampling, 5000 nodes

Παραλληλία

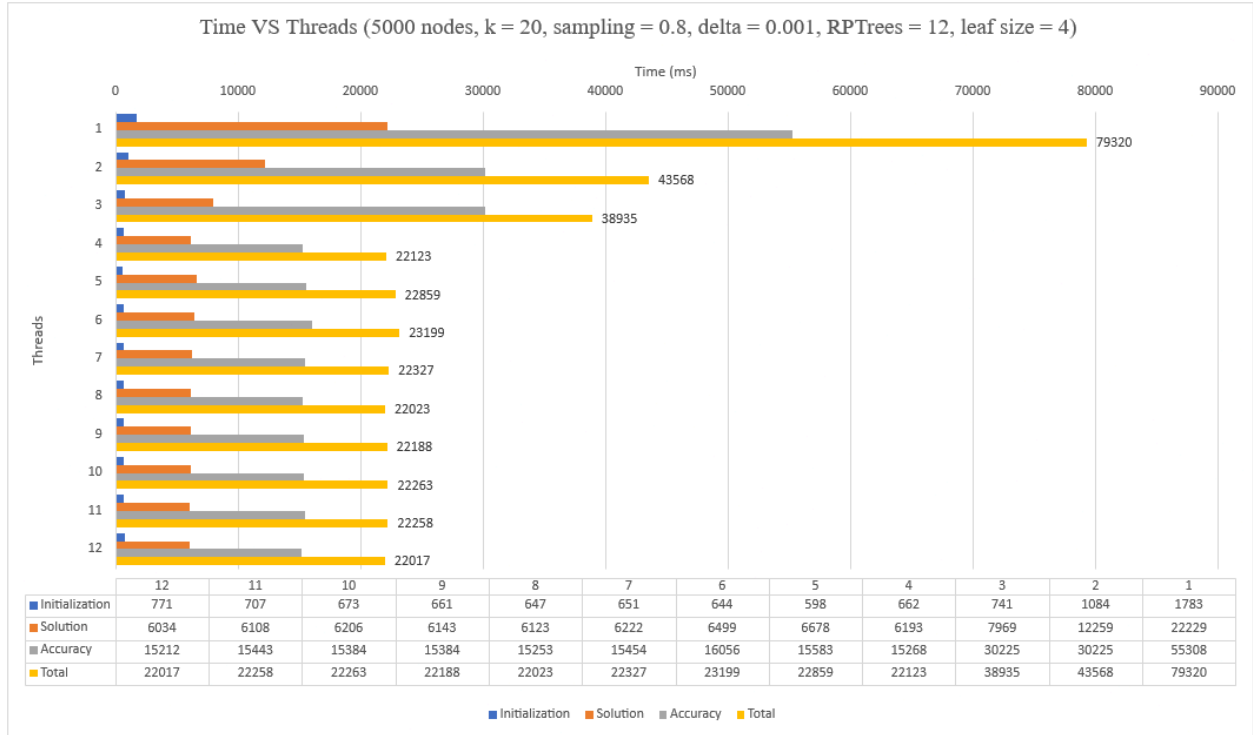
Λόγω της φύσης του προβλήματος (οργάνωση του γράφου σε πίνακα & ανάγκη για επαναληπτικές διαδικασίες), αρκετά στάδια του προβλήματος μπορούσαν να παραλληλοποιηθούν. Η πολυνημάτωση της εφαρμογής μας υλοποιείται μέσω της χρήσης OpenMP. Όλες οι παράλληλες περιοχές μας είναι πάνω σε for loops. Μέσω της χρήσης omp_locks, επιτυγχάνεται το mutual exclusion όταν αυτό χρειάζεται (κυρίως κατα την εισαγωγή γειτόνων).

Τα σημεία όπου παραλληλοποιήθηκαν είναι αρκετά. Ωστόσο, συνοπτικά μπορούμε να αναφέρουμε πως παραλληλοποιήσαμε τα παρακάτω μέρη του κώδικά μας. Όταν αναφερόμαστε σε jobs, αναφερόμαστε στο είδος της διαίρεσης που έγινε σε μια εργασία, π.χ. για τη δημιουργία x δέντρων τυχαίων προβολών, η δουλειά που θα έχει ένα νήμα κάθε φορά θα είναι η δημιουργία ενός δέντρου (job = 1 tree).

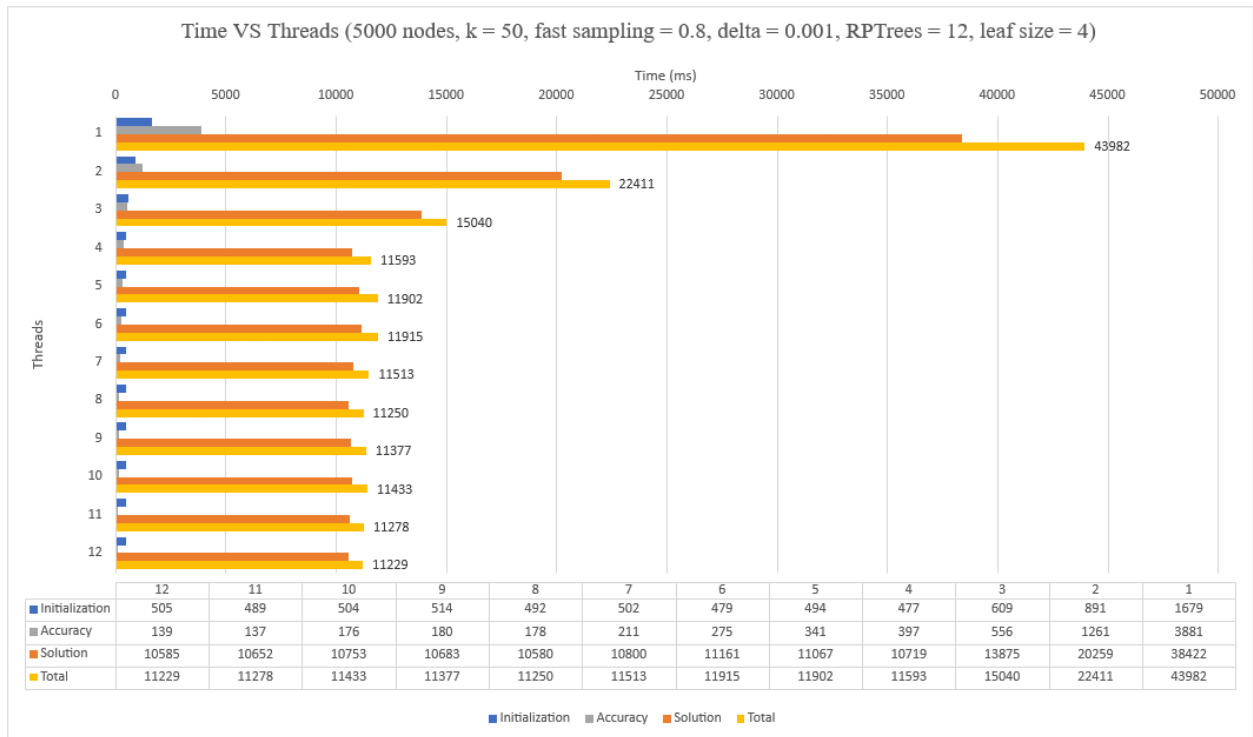
- Δημιουργία δέντρων τυχαίων προβολών (job = 1 tree)
- Αρχικοποίηση γειτόνων μέσω των φύλλων από τα RPTrees (job = 1 leaf)
- Στάδιο 2A: Δειγματοληψία και εύρεση υπονήφειων γειτόνων για κάθε γείτονα (job = 1 neighbor)
- Στάδιο 2B: Ενημέρωση των γειτόνων του KNN γράφου (job = 1 neighbor)
- Υπολογισμός ευστοχίας από αρχείο (παράλληλη ανάγνωση γραμμών από το αρχείο μέσω πολλαπλών file descriptors και χρήση fseek)
- Υπολογισμός νόρμας σημείου
- Υπολογισμός του πραγματικού KNN γράφου

Στο σύστημα στο οποίο κάνουμε τις δοκιμές (οι προδιαγραφές του αναγράφονται παρακάτω), παρατηρούμε ότι με 4 threads φτάνουμε γρήγορα στη μέγιστη αποδοτικότητα. Σίγουρα, η βελτίωση που έρχεται με τη πολυνημάτωση είναι μεγάλη σε μια εφαρμογή όπως αυτή, ειδικά για τον υπολογισμό του πραγματικού γράφου.

Στο εκτελέσιμο **KNN** αλλά και **GraphSolve** μπορούν να δοθούν ως όρισμα ο αριθμός των νημάτων που θα χρησιμοποιούνται από την εκάστοτε εφαρμογή (μέσω της σημαίας -t).



Χρόνος VS νήματα για εκτέλεση του **KNN** (χωρίς λυμένο γράφο), 5000 nodes, k = 20



Χρόνος VS νήματα για εκτέλεση του **KNN** (με λυμένο γράφο), 5000 nodes, k = 50, fast sampling

Δειγματοληψία & Σταθερά δ

Η σταθερά δειγματοληψίας δεν είχε μεγάλο αντίκτυπο στα αποτελέσματα για μικρά k . Οι διακυμάνσεις είναι έως 2% στην ευστοχία, με αντίστοιχες τιμές στην χρονική απόδοση. Ωστόσο, για μεγαλύτερα k , η σταθερά δειγματοληψίας μαζί με την σταθερά δ παίζουν καθοριστικό ρόλο για τον τερματισμό του προγράμματος σε εύλογο χρονικό διάστημα. Ο συνδυασμός μικρό sampling rate/μεγάλο δ για μεγαλύτερα k μπορεί να φέρει πολύ καλύτερο χρόνο με πολύ καλή ευστοχία, ειδικά με την χρήση RPT δέντρων. Ο ιδανικός συνδυασμός εξαρτάται από το k , αλλά και το μέγεθος του dataset.

```
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00010000-4.bin -k 100 -s 0.8 -d 0.4 -r 12 -f
File: 00010000-4.bin | k: 100 | Delta: 0.4 | Sampling rate: 80% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 20 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 2801ms
Solving graph... graph solved at 2 iterations | Time elapsed: 117795ms
Calculating accuracy... relative accuracy: 99.9219% | Absolute accuracy: 99.8417% | Time elapsed: 121621ms
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00010000-4.bin -k 100 -s 0.2 -d 0.4 -r 12 -f
File: 00010000-4.bin | k: 100 | Delta: 0.4 | Sampling rate: 20% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 20 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 1412ms
Solving graph... graph solved at 3 iterations | Time elapsed: 15746ms
Calculating accuracy... relative accuracy: 99.1386% | Absolute accuracy: 98.2657% | Time elapsed: 90245ms
```

VS (fast) sampling rate, 10000 nodes, $k = 100$, $\delta = 0.4$

```
File: 00010000-4.bin | k: 100 | Delta: 0.4 | Sampling rate: 20% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 20 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 1412ms
Solving graph... graph solved at 3 iterations | Time elapsed: 15746ms
Calculating accuracy... relative accuracy: 99.1386% | Absolute accuracy: 98.2657% | Time elapsed: 90245ms
linux18:/home/users/sdi1900156/Documents/Project2023>./KNN 00010000-4.bin -k 100 -s 0.2 -r 12 -f
File: 00010000-4.bin | k: 100 | Delta: 0.001 | Sampling rate: 20% | Fast sampling: 1 | RPTrees: 12 | RPT leaf size: 20 | Threads: 4
Initializing graph... initialization complete | Time elapsed: 2852ms
Solving graph... graph solved at 11 iterations | Time elapsed: 70558ms
Calculating accuracy... relative accuracy: 99.8156% | Absolute accuracy: 99.6256% | Time elapsed: 62714ms
```

VS delta, 10000 nodes, $k = 100$, sampling rate = 20%

Προδιαγραφές συστήματος

Τα προγράμματα μεταγλωττίστηκαν και εκτελέστηκαν σε σύστημα με τα εξής specifications:

- **System**
 - Kernel: 5.4.0-167-generic x86_64
 - Bits: 64
 - Version: 9.4.0
 - Console: tty 6
 - Distro: Ubuntu 20.04.6 LTS (Focal Fossa)
- **CPU**
 - Model: Intel(R) Core(TM) i5-6500 CPU @ 3.20GHz
 - Architecture: x86_64
 - Byte Order: Little Endian
 - Address sizes: 39 bits physical, 48 bits virtual
 - CPU(s): 4
 - Thread(s) per core: 1
 - Core(s) per socket: 4
 - CPU max MHz: 3600.0000
 - CPU min MHz: 800.0000
 - Virtualization: VT-x
 - L1d cache: 128 KiB
 - L1i cache: 128 KiB
 - L2 cache: 1 MiB
 - L3 cache: 6 MiB
- **Compiler:** g++ (Ubuntu 9.4.0-1ubuntu1~20.04.2) 9.4.0