

Τμήμα Μηχανικών Η/Υ και Πληροφορικής

Πανεπιστήμιο Πατρών

23 Νοεμβρίου 2020

Παράλληλος Προγραμματισμός 1^ο Εργαστήριο

Ανδρέας Καρατζάς



Περιεχόμενα

5

ΚΕΦΑΛΑΙΟ 1

Ο αλγόριθμος K-Means

- 1.1 Εισαγωγή 5
- 1.2 Ρυθμίζοντας το Visual Studio 2019 5
- 1.3 Η δομή της εργασίας 7
 - 1.3.1 Ο φάκελος data 7
 - 1.3.2 `Common.h` 7
 - 1.3.3 `Vec.h` και `Vec.cpp` 8
 - 1.3.4 `Center.h` και `Center.cpp` 8
 - 1.3.5 `Class.h` και `Class.cpp` 8
 - 1.3.6 `Optimize.h` και `Optimize.cpp` 8
 - 1.3.7 `Distance.h` και `Distance.cpp` 8
 - 1.3.8 `Interface.h` και `Interface.cpp` 9
 - 1.3.9 `Benchmark.h` και `Benchmark.cpp` 10
 - 1.3.10 `Validation.h` και `Validation.cpp` 10
 - 1.3.11 `kmeans.h` και `kmeans.cpp` 11
- 1.4 Από το testing, στο scaling 11
- 1.5 Τα αποτελέσματα του Intel VTUNE Profiler 12

15

ΚΕΦΑΛΑΙΟ 2

Το πρόγραμμα επικύρωσης σε Python

- 2.1 Περιγραφή 15

17

ΚΕΦΑΛΑΙΟ 3

Ο κώδικας C++

- 3.1 Σχολιασμός 17
- 3.2 Τα αρχεία 17

Ο αλγόριθμος K-Means

1.1 Εισαγωγή

Για την 1^η εργαστηριακή άσκηση, υλοποιήθηκε ο αλγόριθμος K-Means. Ο K-Means είναι ένας αλγόριθμος που χρησιμοποιείται για συσταδοποίηση στο χώρο του Big Data. ο λόγος είναι ότι έχει πολύ καλή απόδοση. Η πολυπλοκότητά του είναι $O(iterations \cdot dataset \cdot K \cdot dimensions)$. Ο ψευδοκώδικας που θα μπορούσε να περιγράψει τον K-Means είναι ο Αλγόριθμος 1.1.

Αλγόριθμος 1.1 K-Means

```

1 input: Vec[N][Nv], Center[Nc][Nv], Classes[N], double threshold
2 output: Center
3 begin
4     while  $\sum_{i=0}^N \sum_{j=0}^{Nc} euclidean\ distance(Vec[i], Center[j]) > threshold$ :
5         Classes  $\leftarrow j \mid \Rightarrow \argmin euclidean\ distance(Vec[i], Center[j]) \forall i \in N \ \& \ j \in Nc$ 
6         foreach Center[i]:
7             Center[i]  $\leftarrow Vec[j] \ \forall j \mid \Rightarrow Vec[j] \in Classes[i]$ 
8     return Center
9 end

```

Στα πλαίσια της εργασίας έγιναν αρκετές βελτιστοποιήσεις. Ωστόσο, ο αλγόριθμος σχεδιάστηκε σε ένα γενικό πλαίσιο και χωρίς τοπικότητα, δηλαδή δεν έγιναν παραδοχές. Για τον έλεγχο του αλγορίθμου οπτικά, φτιάχτηκε και ένα πρόγραμμα σε γλώσσα Python, το οποίο θα αναλυθεί σε επόμενο κεφάλαιο. Στα επόμενα κεφάλαια, θα αναλυθεί η δομή της εργασίας και θα απαντηθούν τα ζητούμενα που τίθενται στην εκφώνηση. Η εργασία υλοποιήθηκε σε C++ 17 με σκοπό την δημοσίευση του κώδικα στην πλατφόρμα του GitHub χρησιμοποιώντας τον IDE του Visual Studio 2019 και το μεταγλωττιστή Intel 19.1. Το σύστημα στο οποίο υλοποιήθηκε ο αλγόριθμος έχει λειτουργικό Windows 10, επεξεργαστή Intel Core i7 9^{ης} γενιάς και μνήμη RAM των 16 GB. Για την εκτέλεση του κώδικα σε λειτουργικό Unix, συντάχθηκε ένα αρχείο makefile το οποίο βρίσκεται στο ίδιο directory με τα αρχεία κώδικα.

1.2 Ρυθμίζοντας το Visual Studio 2019

Αρχικά, σε περίπτωση που το μηχάνημα στο οποίο θα ελεγχθεί η εργασία διαθέτει επεξεργαστή Intel, προτείνεται η εγκατάσταση του Compiler της Intel, καθώς παρέχει επιπλέον βελτιστοποίηση κώδικα για επεξεργαστές της αρχιτεκτονικής τους.

Έπειτα, θα πρέπει να γίνει cloning του repository από το [GitHub](#). Εναλλακτικά, μπορούν να προστεθούν χειροκίνητα τα αρχεία που βρίσκονται στο συμπιεσμένο. Στο συμπιεσμένο υπάρχει και ένα αρχείο *makefile* που μπορεί να αξιοποιηθεί ιδιαιτέρως σε περιβάλλον Unix. Στο περιβάλλον Visual Studio 2019, θα πρέπει να γίνουν αλλαγές στις ρυθμίσεις του Project επιλέγοντας Project > Properties και ακολουθώντας τα παρακάτω βήματα:

- Configuration Properties > General > Platform Toolset → Intel C++ Compiler 19.1
- Configuration Properties > General > C++ Language Standard → ISO C++ 17 Standard (std::c++17)
- Configuration Properties > C++ > Optimization > Optimization → Highest Optimizations (/O3)
- Configuration Properties > C++ > Optimization > Inline Function Expansion → Any Suitable (/Ob2)
- Configuration Properties > C++ > Optimization > Enable Intrinsic Functions → Yes (/Oi)
- Configuration Properties > C++ > Optimization > Favor Size Or Speed → Favor fast code (/Ot)
- Configuration Properties > C++ > Optimization > Enable Fiber-Safe Optimizations → Yes (/GT)
- Configuration Properties > C++ > Optimization [Intel C++] > Optimize For Windows Application → Yes (/GA)
- Configuration Properties > C++ > Code Generation > Enable String Pooling → Yes (/GF)
- Configuration Properties > C++ > Code Generation > Enable C++ Exceptions → Yes (/EHsc)
- Configuration Properties > C++ > Code Generation > Floating Point Model → Fast=2 (/fp:fast=2)
- Configuration Properties > C++ > Code Generation > Enable Floating Point Exceptions → No
- Configuration Properties > C++ > Code Generation [Intel C++] > Floating Point Expression Evaluation → Double (/fp:double)
- Configuration Properties > C++ > Code Generation [Intel C++] > Intel Processor-Specific Optimization → Same as the host processor performing the compilation (/QxHost)

- Configuration Properties > C++ > Language > C++ Language Standard → ISO C++ 17 Standard (/std:c++17)
- Configuration Properties > C++ > Language[Intel C++] > C/C++ Language Support → C++ 17 Support (/Qstd=c++17)

1.3 Η δομή της εργασίας

Ο κώδικας χωρίστηκε σε header files. Στα διάφορα header files δηλώθηκαν συναρτήσεις και μεταβλητές, που ορίζονται στα αντίστοιχα CPP αρχεία κώδικα. Το repository του αλγορίθμου βρίσκεται σε αυτό το [link](#). Στο repository υπάρχει ο φάκελος k-means στον οποίο υπάρχουν όλα τα αρχεία που γράφηκαν στα πλαίσια της εργασίας, μαζί με τα αρχεία `k-means.vcxproj` και `k-means.vcxproj.filters`, τα οποία είναι αρχεία του Visual Studio.

1.3.1 Ο φάκελος data

Στον φάκελο data αποθηκεύονται τα δεδομένα όπως αυτά υπολογίστηκαν από τον αλγόριθμο K-Means. Τα δεδομένα αυτά αποτυπώνουν το τυχαίο dataset, αλλά και τις αλλαγές που υπάρχουν στις μεταβλητές Classes και Centers. Ο χρήστης θα πρέπει να μετακινεί ό,τι αρχεία βρίσκονται μέσα στο συγκεκριμένο φάκελο. Αν και έγινε προσπάθεια σύνταξης κώδικα για να επιλυθεί το παραπάνω πρόβλημα χρησιμοποιώντας τη βιβλιοθήκη [filesystem](#), ο compiler δεν κατάφερε να παράξει εκτελέσιμο.

1.3.2 Common.h

Στην κεφαλίδα `Common.h` δηλώνονται οι απαραίτητες βιβλιοθήκες για την εκτέλεση του προγράμματος, αλλά και σταθερές που καθορίζουν τη λειτουργία του αλγορίθμου. Αναλύοντας:

- **VERBOSITY**: Η σταθερά χρησιμεύει κατά την κλήση της συνάρτησης `kmeans_progress` που είναι δηλωμένη στην κεφαλίδα `Interface`. Η σταθερά αυτή μπορεί να πάρει τιμή 0, 1, 2 ή 3. Όσο πιο μεγάλη τιμή έχει, τόσο μεγαλύτερη θα είναι και η συχνότητα εκτύπωσης των αποτελεσμάτων του `kmeans` κατά τη διαδικασία της βελτιστοποίησης των συστάδων
- **TEST_MODE**: Η σταθερά αυτή χρησιμεύει κατά τη διαδικασία του testing του προγράμματος. Όταν αρχικοποιείται με 1, το πρόγραμμα χρησιμοποιεί σύνολο δεδομένο με πληθάρημο 30 αντί για 100,000 που ορίζεται μέσα στην εκφώνηση. Επίσης, οι διαστάσεις είναι 2 και τα clusters 3. Όταν είναι αρχικοποιημένη με 0, τότε το πρόγραμμα λειτουργεί με τα δεδομένα της εκφώνησης

- `MAX_LIMIT`: Η σταθερά αυτή καθορίζει το μέγιστο που μπορεί να πάρει οποιοδήποτε από τα στοιχεία του συνόλου δεδομένων μας (μεταβλητή `Vec`) έτσι ώστε να λυθεί το πρόβλημα της υπερχείλισης δεδομένων τύπου `float`. Σε περίπτωση που γίνεται `testing`, η σταθερά παίρνει χαμηλότερη τιμή έτσι ώστε να διευκολύνει τον προγραμματιστή

1.3.3 `Vec.h` και `Vec.cpp`

Στο αρχείο `Vec.h` γίνεται η δήλωση της συνάρτησης `vec_init` η οποία καλείται στην αρχή του προγράμματός μας για να γεμίσει τη μεταβλητή `Vec` που αποτελεί το `dataset` του αλγορίθμου. Στο αρχείο `Vec.cpp`, η παραπάνω συνάρτηση εμπλουτίζεται με τον απαραίτητο κώδικα.

1.3.4 `Center.h` και `Center.cpp`

Στο αρχείο `Center.h` γίνεται η δήλωση της συνάρτησης `init_centers` η οποία καλείται στην αρχή του προγράμματός μας για να αρχικοποιηθούν τα `clusters` (μεταβλητή `Center`¹). Στο αρχείο `Center.cpp`, η παραπάνω συνάρτηση εμπλουτίζεται με τον απαραίτητο κώδικα.

1.3.5 `Class.h` και `Class.cpp`

Στο αρχείο `Class.h` γίνεται η δήλωση της συνάρτησης `compute_classes` η οποία καλείται κάθε φορά που πρέπει να αντιστοιχιστούν τα δεδομένα του αλγορίθμου (`dataset`) με το σωστό `cluster` (μεταβλητή `old_Center`). Στο αρχείο `Class.cpp`, η παραπάνω συνάρτηση εμπλουτίζεται με τον απαραίτητο κώδικα.

1.3.6 `Optimize.h` και `Optimize.cpp`

Στο αρχείο `Optimize.h` γίνεται η δήλωση της συνάρτησης `optimize_center` η οποία καλείται σε κάθε επανάληψη του αλγορίθμου και επιστρέφει τα καινούργια διανύσματα - συστάδες (μεταβλητή `new_Center`). Η συνάρτηση αυτή καλείται μετά την αντιστοίχιση δεδομένων με συστάδες, για να βελτιστοποιηθούν τα διανύσματα των συστάδων. Στο αρχείο `Optimize.cpp`, η παραπάνω συνάρτηση εμπλουτίζεται με τον απαραίτητο κώδικα.

1.3.7 `Distance.h` και `Distance.cpp`

Στο αρχείο `Distance.h` γίνεται η δήλωση των παρακάτω συναρτήσεων:

- `euc1_diff`: Η συνάρτηση αυτή δέχεται στην είσοδό της 2 διανύσματα διαστάσεων N_v και επιστρέφει την Ευκλείδεια απόστασή τους

¹Στο πρόγραμμα αναφέρεται ως `old_Center`

- **convergence:** Η συνάρτηση αυτή δέχεται στην είσοδό της τις συστάδες από δύο διαδοχικές επαναλήψεις και επιστρέφει την αθροιστική Ευκλείδεια απόστασή τους. Η συνάρτηση καλείται για να αποφανθεί ο αλγόριθμος σύγκλισης
- **normalize_convergence:** Η συνάρτηση αυτή δέχεται την προηγουμένως υπολογισμένη αθροιστική Ευκλείδεια απόσταση δύο διαδοχικών συστάδων και την κανονικοποιεί σύμφωνα με το μαθηματικό τύπο:

$$\frac{|A - B|}{\max A, B}$$

Ωστόσο, αν ο αλγόριθμος είναι ακόμα στις πρώτες επαναλήψεις, όπου η Ευκλείδεια απόσταση αλλάζει ραγδαία, τότε η κανονικοποιημένη τιμή θεωρείται ίση με 1 (μέγιστη κανονικοποιημένη τιμή). Τέλος, για να λυθεί το πρόβλημα του υπολογιστικού σφάλματος λόγω έλλειψης επαρκούς ακρίβειας σε αριθμούς κινητής υποδιαστολής, εισάγεται ένας επιπλέον κανόνας: Αν η μη κανονικοποιημένη αθροιστική Ευκλείδεια απόσταση δύο διαδοχικών συνόλων συστάδων είναι μικρότερη του 100, τότε θεωρείται ότι η επιτάχυνση των κέντρων είναι σχεδόν μηδενική² και άρα ο αλγόριθμος έχει συγκλίνει

Στο αρχείο `Distance.cpp`, οι συναρτήσεις αυτές εμπλουτίζονται με τον αντίστοιχο κώδικα.

1.3.8 Interface.h και Interface.cpp

Στο αρχείο `Interface.h` γίνεται η δήλωση των παρακάτω συναρτήσεων:

- **multidimensional_float_vector_interface:** Η συνάρτηση αυτή εκτυπώνει σε περιβάλλον CLI περιεχόμενα οποιασδήποτε μεταβλητής τύπου `std::vector<std::array<float, Nv>>`. Στο πρόγραμμα, χρησιμοποιείται για την εκτύπωση των περιεχομένων των μεταβλητών:
 - `Vec`,
 - `old_Center` και
 - `new_Center`
- **multidimensional_int_array_interface:** Η συνάρτηση αυτή εκτυπώνει σε περιβάλλον CLI περιεχόμενα οποιασδήποτε μεταβλητής τύπου `std::array<std::vector<int>, Nc>`. Στο πρόγραμμα, χρησιμοποιείται για την εκτύπωση των περιεχομένων της μεταβλητής `Classes`
- **progress_interface και kmeans_progress:** Η συνάρτηση `progress_interface` καλείται από τη συνάρτηση `kmeans_progress`, η οποία πληροφορεί το χρήστη για τη σχετική πρόοδο του αλγορίθμου

²Η επιτάχυνση των κέντρων είναι πολύ μικρή σε αυτήν την περίπτωση καθώς η μέση μετατόπιση σημείου μεταξύ των επαναλήψεων αυτών είναι μικρότερη του $\frac{100}{100 \cdot 1,000} = 0.001$.

- `kmeans_termination`: Η συνάρτηση αυτή καλείται στο τέλος του αλγορίθμου και εκτυπώνει τους χρόνους (benchmark) του αλγορίθμου και επιβεβαιώνει τη σύγκλισή του

Στο αρχείο `Interface.cpp`, οι συναρτήσεις αυτές εμπλουτίζονται με τον αντίστοιχο κώδικα.

1.3.9 `Benchmark.h` και `Benchmark.cpp`

Στο αρχείο `Benchmark.h` δηλώνονται οι συναρτήσεις:

- `init_benchmark`: Η συνάρτηση αυτή επιστρέφει ένα σημείο στο χρόνο (checkpoint), οπουδήποτε το χρειάζεται ο προγραμματιστής, για να μετρήσει χρόνο (benchmarking)
- `bench_convergence`: Η συνάρτηση αυτή θα μπορούσε να διαγραφεί, καθώς έχει ίδιο σώμα με τη συνάρτηση `init_benchmark`. Ωστόσο, λόγω της διαφοράς στο όνομα, βοηθάει τον προγραμματιστή να ξεχωρίζει το σκοπό της καθεμίας και συνεπώς το που πρέπει να κληθεί η καθεμία
- `bench_loop`: Η συνάρτηση αυτή μετράει το συνολικό χρόνο που βρίσκεται «εγκλωβισμένος» ο αλγόριθμος μέσα στην επανάληψη βελτιστοποίησης
- `terminate_bench`: Η συνάρτηση αυτή θα μπορούσε να διαγραφεί, καθώς έχει ίδιο σώμα με τη συνάρτηση `init_benchmark` και τη `bench_convergence`. Ωστόσο, λόγω της διαφοράς στο όνομα, βοηθάει τον προγραμματιστή να ξεχωρίζει το σκοπό της καθεμίας και συνεπώς το που πρέπει να κληθεί η καθεμία
- `benchmark_results`: Η συνάρτηση αυτή χρησιμοποιείται στο τέλος του προγράμματος κι επιτελεί συνολικό benchmarking του αλγορίθμου

Στο αρχείο `Benchmark.cpp`, οι συναρτήσεις αυτές εμπλουτίζονται με τον αντίστοιχο κώδικα.

1.3.10 `Validation.h` και `Validation.cpp`

Στο αρχείο `Validation.h` δηλώνονται οι συναρτήσεις:

- `export_multidimensional_float_vector`: Η συνάρτηση αυτή αποθηκεύει σε αρχείο CSV τα περιεχόμενα οποιασδήποτε μεταβλητής τύπου `std::vector<std::array<float, N>>`. Στο πρόγραμμα, χρησιμοποιείται για την εκτύπωση των περιεχομένων των μεταβλητών:
 - `Vec`,
 - `old_Center` και
 - `new_Center`

- `export_multidimensional_integer_array`: Η συνάρτηση αυτή αποθηκεύει σε αρχείο CSV τα περιεχόμενα οποιασδήποτε μεταβλητής τύπου `std::array<std::vector<int>, Nc>`. Στο πρόγραμμα, χρησιμοποιείται για την εκτύπωση των περιεχομένων της μεταβλητής `Classes`
- `track_kmeans_progress` και `export_kmeans_progress`: Η συνάρτηση `track_kmeans_progress` ενημερώνει τη μεταβλητή `kmeans_progress` την οποία αποθηκεύει σε αρχείο TXT η συνάρτηση `export_kmeans_progress`. Τα δεδομένα που εξάγονται αφορούν τη πρόοδο που κάνει ο αλγόριθμος K-Means από επανάληψη σε επανάληψη

Ο σκοπός αυτών των συναρτήσεων είναι εμφανής κατά το testing του προγράμματος. Τα δεδομένα που εξήχθησαν στο φάκελο `data` μετά από την εκτέλεση του αλγορίθμου, τα επεξεργάζεται το πρόγραμμα σε Python, το οποίο επιστρέφει τη γεωμετρική εικόνα της πορείας του αλγορίθμου. Στο αρχείο `Validation.cpp`, οι συναρτήσεις αυτές εμπλουτίζονται με τον αντίστοιχο κώδικα.

1.3.11 `kmeans.h` και `kmeans.cpp`

Στο αρχείο `kmeans.cpp` δηλώνονται όλες οι επιπλέον (custom) κεφαλίδες του προγράμματος. Στο αρχείο `kmeans.cpp` είναι γραμμένος ο driver του προγράμματος (συνάρτηση `main`).

1.4 Από το testing, στο scaling

Το πρόγραμμα όπως έχει ανέβει στο GitHub είναι αρχικοποιημένο για testing. Για να μπορέσει ο εκάστοτε χρήστης να δει την καθαρή απόδοση του προγράμματος³, θα πρέπει να γίνει η ακόλουθη επεξεργασία:

- `Common.h` > line 27 > `TEST_MODE ← 0`
- `kmeans.cpp` > lines [18, 19] ← Comment
- `kmeans.cpp` > lines [23, 24] ← Comment
- `kmeans.cpp` > line 36 ← Comment
- `kmeans.cpp` > lines [38, 39] ← Comment
- `kmeans.cpp` > lines [41, 42] ← Comment
- `kmeans.cpp` > line 47 ← Comment
- `kmeans.cpp` > lines [52, 53] ← Comment

³Το πρόγραμμα φαίνεται αν είναι αποδοτικό ή όχι, αν οι μεταβλητές `N`, `Nv` και `Nc` αρχικοποιηθούν όπως περιγράφεται στο τέταρτο κομμάτι της εργασίας.

Σε περίπτωση που ο tester επιθυμεί να κάνει profiling ή να χρησιμοποιήσει την `time` λειτουργικού Unix, τότε θα πρέπει να κάνει επιπλέον των παραπάνω και τις παρακάτω αλλαγές:

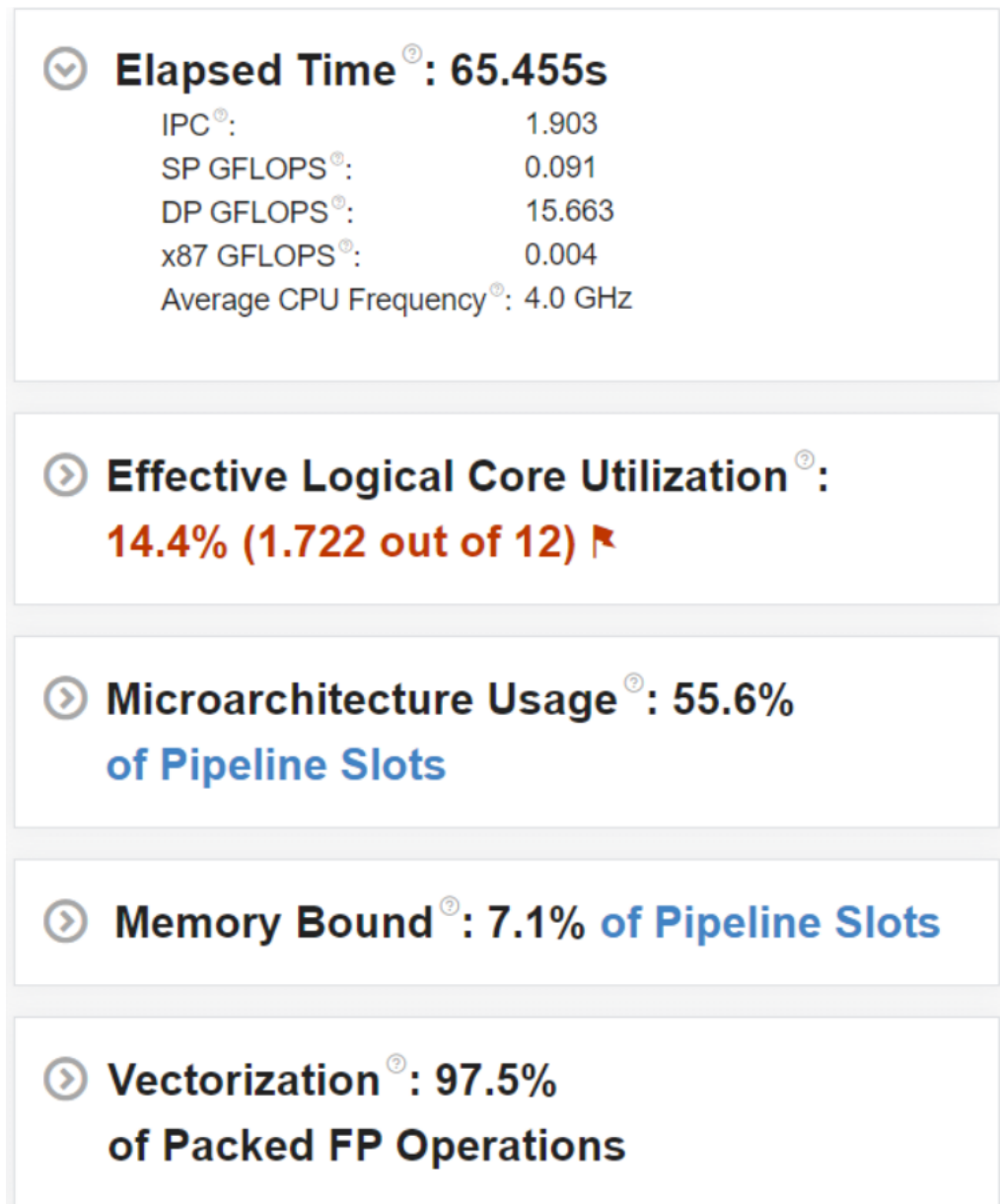
- `kmeans.cpp` > lines [12, 14] ← Comment
- `kmeans.cpp` > line 30 ← Comment
- `kmeans.cpp` > lines [33, 35] ← Comment
- `kmeans.cpp` > lines [49, 51] ← Comment

1.5 Τα αποτελέσματα του Intel VTUNE Profiler

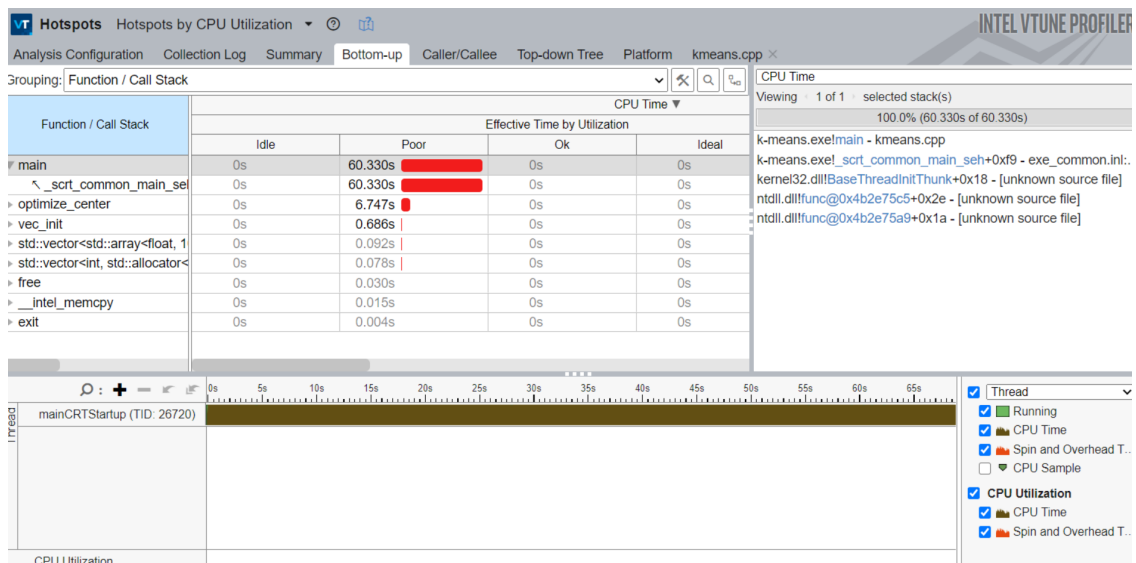
Το αποτέλεσμα του Intel VTUNE Profiler φαίνεται στο Σχήμα 1.1.

Το πρώτο στοιχείο που φαίνεται είναι ο χρόνος που χρειάστηκε το πρόγραμμα για την εκτέλεσή του. Ο χρόνος είναι μέσα στα επιθυμητά πλαίσια όπως αυτά ορίζονται στην εκφώνηση. Η κόκκινη ένδειξη δείχνει το πόση παραλληλοποίηση έγινε στον αλγόριθμο. Επειδή η άσκηση δε ζητούσε παραλληλοποίηση, το πρόγραμμα είναι 1 - *threaded*. Θετικά είναι τα ποσοστά *Memory Bound* που δείχνει την αποδοτικότητα των I/Os που έγιναν κατά την εκτέλεση του προγράμματος, και το *Vectorization*, το οποίο επιτρέπει block επεξεργασία.

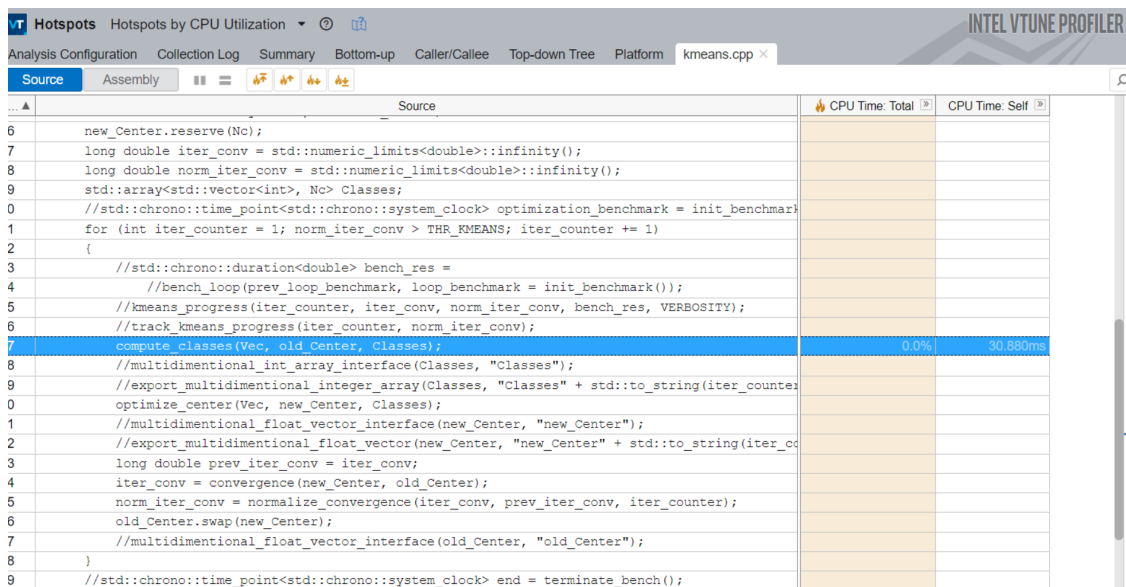
Μια λεπτομερέστερη εικόνα του προγράμματος δίνεται στα σχήματα 1.2 και 1.3, όπου φαίνεται ανά συνάρτηση η κατανάλωση χρόνου. Βλέπουμε ότι ο Profiler δίνει την `main` που είναι λογικό, καθώς αποτελεί το `driver` του προγράμματος. Ωστόσο, στο Σχήμα 1.3 ξεκαθαρίζει λίγο παραπάνω η κατανομή και πλέον δίνεται η συνάρτηση `compute_classes`, η οποία καταναλώνει σχεδόν το μισό χρόνο του προγράμματος. Αυτό είναι απολύτως λογικό, καθώς καλείται σε κάθε επανάληψη του K-Means και έχει περισσότερες πράξεις από κάθε άλλη συνάρτηση εντός της επανάληψης.



Σχήμα 1.1: Η ανάλυση του Intel VTUNE Profiler



Σχήμα 1.2: Η ανάλυση της main



Σχήμα 1.3: Το «hotspot» του προγράμματος

Το πρόγραμμα επικύρωσης σε Python

2.1 Περιγραφή

Για τον ευκολότερο έλεγχο του προγράμματος συντάχθηκε και ένα πρόγραμμα σε Python με περιβάλλον *PyCharm*. Το πρόγραμμα υπάρχει στο φάκελο `kmeans-visualize`. Το πρόγραμμα χρειάζεται τις βιβλιοθήκες:

- [NumPy](#)
- [Pandas](#)
- [Matplotlib](#)

Προαιρετικά, ο προγραμματιστής μπορεί να εγκαταστήσει και το πακέτο [SciPy](#). Έπειτα, ο προγραμματιστής μετακινεί τα αρχεία επικύρωσης που δημιουργήθηκαν στο φάκελο `data` από τον αλγόριθμο K-Means στο φάκελο που υπάρχει και το Python Script. Τέλος, ο προγραμματιστής εκτελεί το Python Script.

Στο GitHub υπάρχουν και [οδηγίες](#) για το πως μπορεί ο προγραμματιστής να μετατρέψει αυτό το «stream» αρχείων PNG που δημιουργούνται μετά από την εκτέλεση του Script σε αρχείο GIF (βλέπε αρχείο [«instructions»](#)).

Ο κώδικας C++

3.1 Σχολιασμός

Στην αρχή η ιδέα ήταν να χρησιμοποιηθούν αποκλειστικά πίνακες για την αποθήκευση όλων των δεδομένων. Ωστόσο, αυτό αποδείχθηκε εξαιρετικά απαιτητικό ως προς το stack του προγράμματος. Έτσι, έγινε η αλλαγή σε δυναμικούς πίνακες που αποτελούνται από (στατικούς) πίνακες. Το αποτέλεσμα ήταν ότι κάθε φορά ήταν επαρκές να υπάρχει στη μνήμη ένα διάνυσμα 1000 στοιχείων, που θα αποτελούσε το διάνυσμα προς επεξεργασία από τη μεταβλητή `Vec`, ένα διάνυσμα επίσης 1000 στοιχείων που θα αποτελούσε κάποιο από τα κέντρα της μεταβλητής `old_Center` ή `new_Center` (ανάλογα με το σημείο εκτέλεσης) και οι δείκτες προς όλους τους πίνακες οι οποίοι όμως είναι αποθηκευμένοι στο heap. Στο φάκελο `buggy` υπάρχουν ενδεικτικά κάποιες ενδιάμεσες εκδόσεις που ικανοποιούν το ζητούμενο της εργασίας 4. Ωστόσο, δεν υπάρχουν σχόλια στις εκδόσεις αυτές, καθώς είναι κάποιες ειδικά έχουν αρκετά λανθασμένη αρχιτεκτονική.

3.2 Τα αρχεία

Κώδικας 3.2.1: Η κεφαλίδα Common.h

```
1 /**
2  * Common.h
3  *
4  * In this header file, we define the constants
5  * used throughout the K-means algorithm. We also
6  * include all the header files necessary to make
7  * the implementation work.
8  */
9
10 #pragma once
11
12 #include <iostream>           ///< std::cout
13 #include <vector>             ///< std::vector
14 #include <iomanip>            ///< std::setw
15 #include <limits>             ///< std::numeric_limits
16 #include <chrono>            ///< std::chrono
17 #include <array>             ///< std::array
18 #include <algorithm>         ///< std::find
19 #include <utility>           ///< std::pair
20 #include <string>            ///< std::string
21 #include <cstdlib>           ///< std::abs
22 #include <random>            ///< std::random_device
23 #include <fstream>           ///< std::ofstream
24 #include <numeric>           ///< std::iota
25
26 constexpr int VERBOSITY = 3;           ///< Sets K-means' interface verbosity level
27 constexpr int TEST_MODE = 1;          ///< Set it equal to 1 to enter test mode, otherwise set it to 0
28 constexpr int MAX_LIMIT = TEST_MODE == 1 ? 31 : 63; ///< Upper bound used in the randomly generated dataset of reals
29 constexpr int N = TEST_MODE == 1 ? 20 : 100000;   ///< Elements in our dataset
30 constexpr int Nv = TEST_MODE == 1 ? 2 : 1000;     ///< Dimentions of each element
31 constexpr int Nc = TEST_MODE == 1 ? 3 : 100;      ///< Number of clusters, also known as `K`
32 constexpr double THR_KMEANS = 0.000001;          ///< Threshold used to indicate K-means convergence
```

Κώδικας 3.2.2: Η κεφαλίδα Vec.h

```
1 /**
2  * Vec.h
3  *
4  * In this header file, we define a function that
5  * executes only in the beginning of the algorithm.
6  * This function is used to generate a random dataset
7  * for the K-means algorithm.
8  */
9
10 #pragma once
11
12 #include "Common.h"
13
14 void vec_init(std::vector<std::array<float, Nv>>& Vec);
```

Κώδικας 3.2.3: Το αρχείο Vec.cpp

```
1 #include "Vec.h"
2
3 /**
4  * Initializes `Vec` variable.
5  *
6  * @param[in, out] Vec the random dataset generated for clustering using the K-means algorithm.
7  *
8  * @see Common.h
9  *
10  * @remark [<random> Engines and Distributions](https://docs.microsoft.com/en-us/cpp/standard-library/random?view=msvc-160#engdist)
11  */
12 void vec_init(std::vector<std::array<float, Nv>>& Vec)
13 {
14     std::random_device rd_Vec;           ///< Initializes a random generator
15     std::mt19937 mt_Vec(rd_Vec());       ///< Uses the mt19937 engine
16     std::uniform_real_distribution<float> dist_Vec(0.0, MAX_LIMIT + 0.0);      ///< Generates a uniform distribution bounded by the `MAX_LIMIT` set in `Common.h` to prevent number overflow
17     for (int i = 0; i < N; i += 1)
18     {
19         std::array<float, Nv> Elements;   ///< Declares a vector `Elements` that temporarily holds the vector that is to be inserted to `Vuc`
20         for (int j = 0; j < Nv; j += 1)
21         {
22             Elements[j] = dist_Vec(mt_Vec);           ///< Updates contents of `Elements`
23         }
24         Vec.emplace_back(Elements);               ///< Moves `Elements` to `Vec`
25     }
26 }
```

Κώδικας 3.2.4: Η κεφαλίδα Center.h

```
1 /**
2  * Center.h
3  *
4  * In this header file, we define a function that helps
5  * in the random initialization of the variable `Center`.
6  * This function is only used in the beginning of the
7  * K-means algorithm.
8  */
9
10 #pragma once
11
12 #include "Common.h"
13 #include "Distance.h"
14
15 void init_centers(const std::vector<std::array<float, Nv>>& Vec, std::vector<std::array<float, Nv>>& old_Center);
```

Κώδικας 3.2.5: Το αρχείο Center.cpp

```
1 #include "Center.h"
2
3 /**
4  * Initializes `Center` variable.
5  *
6  * @param[in] Vec the dataset generated.
7  * @param[in, out] old_Center the random centroids generated in the beginning of the K-means algorithm.
8  *
9  * @remark [<random> Engines and Distributions](https://docs.microsoft.com/en-us/cpp/standard-library/random?view=msvc-160#engdist)
10  */
11 void init_centers(const std::vector<std::array<float, Nv>>& Vec, std::vector<std::array<float, Nv>>& old_Center)
12 {
13     std::vector<int> Random(N);           ///< Declares our random `Vec` generator vector
14     std::iota(Random.begin(), Random.end(), 0);           ///< Initializes vector with all possible indeces of `Vec`
15     std::shuffle(Random.begin(), Random.end(), std::mt19937{ std::random_device{}() });           ///< Shuffles `Random` vector
16     for (int i = 0; i < Nc; i += 1)
17     {
18         std::array<float, Nv> Element;           ///< Declares a vector `Elements` that temporarily holds the vector chosen from `Vec`
19         Element = Vec.at(Random.at(i));           ///< Copies contents of the random `Vec` element
20         old_Center.emplace_back(Element);           ///< Moves `Element` variable to the `old_Center` variable
21     }
22 }
```

Κώδικας 3.2.6: Η κεφαλίδα Class.h

```
1 /**
2  * Class.h
3  *
4  * In this header file, we a function that assigns
5  * the vectors found in the variable `Vec` to a
6  * corresponding `Class` or more commonly known as cluster.
7  */
8
9 #pragma once
10
11 #include "Common.h"
12 #include "Distance.h"
13
14 void compute_classes(const std::vector<std::array<float, Nv>>& Vec, const std::vector<std::array<float, Nv>>& old_Center, std::array<std::vector<int>, Nc>& Classes);
```

Κώδικας 3.2.7: Το αρχείο Class.cpp

```
1 #include "Class.h"
2
3 /**
4  * Computes `Class` variable.
5  *
6  * @param[in] Vec the dataset generated.
7  * @param[in] old_Center the current centroids.
8  * @param[in, out] Classes the classes corresponding to the dataset.
9  *
10  * @note This function calls `eucl_diff` from Distance.h
11  *
12  * @see [K-means clustering Algorithm](https://medium.com/@jaredchillers_38839/k-means-clustering-algorithm-4334db89bdf3)
13  */
14 void compute_classes(const std::vector<std::array<float, Nv>>& Vec, const std::vector<std::array<float, Nv>>& old_Center, std::array<std::vector<int>, Nc>& Classes)
15 {
16     for (int i = 0; i < Nc; i += 1)                                ///< Clears contents of `Classes` variable
17     {
18         Classes[i].clear();                                         ///< Prevents garbage processing
19         Classes[i].reserve((int) N/Nc);                             ///< Optimizes array with high probability
20     }
21     int argmin_idx = -1;                                           ///< Initializes `argmin` index
22     long double argmin_val = std::numeric_limits<long int>::max() + 0.0; ///< Initializes `argmin` value
23     for (int i = 0; i < N; i += 1)                                  ///< Loop through `Vec`
24     {
25         for (int j = 0; j < Nc; j += 1)                            ///< Loop through `old_Center`
26         {
27             long double temp_eucl_dist = eucl_diff(Vec.at(i), old_Center.at(j)); ///< Compute Euclidean distance between parsed `Vec` element and parsed `old_Center` element
28             if (argmin_val > temp_eucl_dist)                        ///< Checks if this Euclidean distance is so far the smallest
29             {
30                 argmin_val = temp_eucl_dist;                      ///< Updates `argmin` value
31                 argmin_idx = j;                                    ///< Updates `argmin` index
32             }
33         }
34         Classes[argmin_idx].emplace_back(i);                      ///< Assigns `Vec` element to a cluster
35         argmin_idx = -1;                                           ///< Sets `argmin` index for the next loop
36         argmin_val = std::numeric_limits<unsigned long int>::max(); ///< Sets `argmin` value for the next loop
37     }
38 }
```

Κώδικας 3.2.8: Η κεφαλίδα Optimize.h

```
1 /**
2  * Optimize.h
3  *
4  * In this header file, we define a function which
5  * computes the error of the `Center` variable with
6  * respect to the `Vec` variable.
7  */
8
9 #pragma once
10
11 #include "Common.h"
12
13 void optimize_center(const std::vector<std::array<float, Nv>>& Vec, std::vector<std::array<float, Nv>>& new_Center, const std::array<std::vector<int>, Nc>& Classes);
```

Κώδικας 3.2.9: Το αρχείο Optimize.cpp

```
1 #include "Optimize.h"
2
3 /**
4  * Computes `Center` variable after `Classes` correction.
5  *
6  * @param[in] Vec the dataset.
7  * @param[in] Classes the classes corresponding to the dataset.
8  * @param[in, out] new_Center the new clusters of the K-means algorithm.
9  */
10 void optimize_center(const std::vector<std::array<float, Nv>>& Vec, std::vector<std::array<float, Nv>>& new_Center, const std::array<std::vector<int>, Nc>& Classes)
11 {
12     new_Center.clear();                                           ///< Clears previously computed data
13     for (int i = 0; i < Nc; i += 1)                                ///< Loops through all clusters
14     {
15         std::array<float, Nv> Element = { 0 };                    ///< Initializes `Element` array which holds a cluster vector
16         for (int j = 0; j < Classes[i].size(); j += 1)            ///< Loops through `Vec` vectors
17         {
18             for (int k = 0; k < Nv; k += 1)                        ///< Loops through vector elements
19             {
20                 Element[k] += (float) (Vec.at(Classes[i].at(j))[k] / Classes[i].size()); ///< Computes vector element
21             }
22         }
23         new_Center.emplace_back(Element);                          ///< Updates `new_Center` variable
24     }
25 }
```

Κώδικας 3.2.10: Η κεφαλίδα Distance.h

```
1 /**
2  * Distance.h
3  *
4  * In this header file, we define some functions which
5  * help with the computations needed to define whether
6  * K-means has converged or not.
7  */
8
9 #pragma once
10
11 #include "Common.h"
12
13 long double eucl_diff(const std::array<float, Nv>& src, const std::array<float, Nv>& dst);
14 long double convergence(const std::vector<std::array<float, Nv>>& curr_Center, const std::vector<std::array<float, Nv>>& prev_Center);
15 long double normalize_convergence(const long double curr_iter_conv, const long double prev_iter_conv, const int iter_counter);
```



```
Κώδικας 3.2.11: Το αρχείο Distacne.cpp
1  #include "Distance.h"
2
3  /**
4   * Computes Euclidean distance between 2 arrays of size `Nv`.
5   *
6   * @param[in] A the first array.
7   * @param[in] B the second array.
8   *
9   * @return Euclidean distance.
10  */
11  long double eucl_diff(const std::array<float, Nv>& A, const std::array<float, Nv>& B)
12  {
13      long double eucl_diff = 0;
14      for (int i = 0; i < Nv; i += 1)                                     ///< Loops through all `Nv` elements of the given arrays
15      {
16          long double point_diff = A[i] - B[i];
17          if (point_diff > std::numeric_limits<float>::max())             ///< Computes element-wise difference
18          {                                                             ///< Checks for data overflow
19              point_diff = std::numeric_limits<float>::max();
20          }
21          long double point_diff_squared = point_diff * point_diff;
22          eucl_diff += point_diff_squared;                               ///< Squares difference
23      }                                                                ///< Adds difference to `eucl_diff`
24      return eucl_diff;
25  }
26
27  /**
28   * Checks if K-means has converged.
29   *
30   * @param[in] curr_Center the current `Center` variable .
31   * @param[in] prev_Center the old `Center` variable.
32   *
33   * @return Euclidean distance between the 2 vectors.
34   *
35   * @note if the Euclidean distance between 2 points is not greater than 1.0, then it is normalized to 0.
36   *       This is to optimize speed, prevent overfitting, and solve the precision error carried throughout the  $N_c \times N_v = 1,000,000$  additions.
37   */
38  long double convergence(const std::vector<std::array<float, Nv>>& curr_Center, const std::vector<std::array<float, Nv>>& prev_Center)
39  {
40      long double convergence_sum = 0;
41      for (int i = 0; i < Nc; i += 1)                                     ///< Initializes `convergence_sum` variable used to store the total additive difference between `curr_Center` and `prev_Center`
42      {                                                                 ///< Loops through all centers
43          long double tmp_eucl_d = eucl_diff(curr_Center.at(i), prev_Center.at(i));
44          convergence_sum += tmp_eucl_d > 1.0 ? tmp_eucl_d : 0.0;         ///< Computes element-wise Euclidean distance
45          if (convergence_sum > std::numeric_limits<double>::max())       ///< Optimizes Convergence if error is close to defined threshold
46          {                                                             ///< Checks for number overflow
47              convergence_sum = std::numeric_limits<double>::max();
48              break;                                                     ///< Prevents number overflow
49          }                                                             ///< Breaks loop because `convergence_sum` is mazed
50      }
51      return convergence_sum;
52  }
53
54  /**
55   * Normalizes the error between current `Center` variable and old `Center` variable.
56   *
57   * @param[in] curr_iter_conv the current `Center` error.
58   * @param[in] prev_iter_conv the old `Center` error.
59   * @param[in] iter_counter the loop number.
60   *
61   * @return normalized Euclidean distance between centers.
62   *
63   * @note The computation of the normalization formula is explained below:
64   *       if the current loop error is less or equal to the number of clusters, then normalized = 0.0
65   *       if the current iteration of the optimization loop is less or equal to 2, then normalized = 1.0
66   *       else, normalized =  $\frac{|curr\_iter\_conv - prev\_iter\_conv|}{\arg \max curr\_iter\_conv, prev\_iter\_conv}$ 
67   */
68  long double normalize_convergence(const long double curr_iter_conv, const long double prev_iter_conv, const int iter_counter)
69  {
70      return curr_iter_conv > Nc && iter_counter > 2 ? std::abs(curr_iter_conv - prev_iter_conv) / std::max(curr_iter_conv, prev_iter_conv) :
71      iter_counter <= 2 ? 1.0 : 0.0;
72  }
```

```
Κώδικας 3.2.12: Η κεφαλίδα Interface.h
1  /**
2   * interface.h
3   *
4   * In this header file, we define some functions which
5   * implement a basic User Interface (UI) throughout the
6   * execution of the K-means algorithm.
7   */
8
9  #pragma once
10
11  #include "Common.h"
12
13  void multidimentional_float_vector_interface(const std::vector<std::array<float, Nv>>& Obj, std::string Obj_name);
14  void multidimentional_int_array_interface(const std::array<std::vector<int>, Nc>& Obj, std::string Obj_name);
15  void progress_interface(const int iter_counter, const long double iter_conv, const long double norm_iter_conv, const std::chrono::duration<double> loop_benchmark);
16  void kmeans_progress(const int iter_counter, const long double iter_conv, const long double norm_iter_conv, const std::chrono::duration<double> loop_benchmark, const int verbose);
17  void kmeans_termination(const std::pair<std::chrono::duration<double>, std::chrono::duration<double>> bench_results, const long double iter_conv);
```

```

Κώδικας 3.2.13: Το αρχείο Interface.cpp
1  #include "Interface.h"
2
3  /**
4   * Prints the contents of a vector of arrays.
5   *
6   * @param[in] Obj the vector of arrays.
7   * @param[in] obj_name the name of the variable.
8   *
9   * @note This function helps in printing the contents of the variables `Vec`, `old_Centers` and `new_Centers`.
10  */
11 void multidimentional_float_vector_interface(const std::vector<std::array<float, Nv>>& Obj, std::string Obj_name)
12 {
13     for (int i = 0; i < Obj.size(); i += 1)
14     {
15         for (int j = 0; j < Obj.at(i).size(); j += 1)
16         {
17             std::cout << std::setw(11) << Obj_name << "[" << i << "]"[" " << j << "]:\t" << std::setw(18) << Obj.at(i)[j] << " ";
18         }
19         std::cout << std::endl;
20     }
21     std::cout << std::endl << std::endl;
22 }
23
24 /**
25 * Prints the contents of an array of vectors.
26 *
27 * @param[in] Obj the array of vectors.
28 * @param[in] obj_name the name of the variable.
29 *
30 * @note This function helps in printing the contents of the variable `Classes`.
31 */
32 void multidimentional_int_array_interface(const std::array<std::vector<int>, Nc>& Obj, std::string Obj_name)
33 {
34     for (int i = 0; i < Obj.size(); i += 1)
35     {
36         for (int j = 0; j < Obj[i].size(); j += 1)
37         {
38             std::cout << Obj_name << "[" << i << "]"[" " << j << "]: " << std::setw(4) << Obj[i].at(j) << " ";
39             if (j % 6 == 0 && j != 0)
40             {
41                 std::cout << std::endl;
42             }
43         }
44         std::cout << std::endl << std::endl;
45     }
46     std::cout << std::endl << std::endl;
47 }
48
49 /**
50 * Prints the progress of the K-means algorithm.
51 *
52 * @param[in] iter_counter the current loop number.
53 * @param[in] iter_conv the current loop error.
54 * @param[in] norm_iter_conv the current normalized loop error.
55 * @param[in] loop_benchmark the current loop benchmark
56 */
57 void progress_interface(const int iter_counter, const long double iter_conv, const long double norm_iter_conv, const std::chrono::duration<double> loop_benchmark)
58 {
59     std::cout << "Iteration: " << std::setw(2) << iter_counter
60     << "\tError: " << std::setw(std::numeric_limits<float>::max_digits10) << iter_conv
61     << "\tNormalized error: " << std::setw(std::numeric_limits<float>::max_digits10) << norm_iter_conv
62     << "\tLoop benchmark: " << std::setw(std::numeric_limits<double>::max_digits10) << loop_benchmark.count();
63     TEST_MODE == 1 ? (std::cout << std::endl << std::endl << std::endl) : std::cout << std::endl;
64 }
65
66 /**
67 * Calls `progress_interface` function to print K-means progress.
68 *
69 * @param[in] iter_counter the current loop number.
70 * @param[in] norm_iter_conv the current normalized loop error.
71 * @param[in] iter_conv the current loop error.
72 * @param[in] loop_benchmark the current loop benchmark
73 * @param[in] verbose the frequency used for activating the interface
74 * if equal to 0, then `progress_interface` is called every 4 loops
75 * if equal to 1, then `progress_interface` is called every 3 loops
76 * if equal to 2, then `progress_interface` is called every 2 loops
77 * if equal to 3, then `progress_interface` is called in every loop
78 *
79 * @see `progress_interface`
80 */
81 void kmeans_progress(const int iter_counter, const long double iter_conv, const long double norm_iter_conv, const std::chrono::duration<double> loop_benchmark, const int verbose)
82 {
83     switch (verbose)
84     {
85     case 0: if (iter_counter % 4 == 0) { progress_interface(iter_counter, iter_conv, norm_iter_conv, loop_benchmark); }
86             break;
87     case 1: if (iter_counter % 3 == 0) { progress_interface(iter_counter, iter_conv, norm_iter_conv, loop_benchmark); }
88             break;
89     case 2: if (iter_counter % 2 == 0) { progress_interface(iter_counter, iter_conv, norm_iter_conv, loop_benchmark); }
90             break;
91     case 3: if (iter_counter % 1 == 0) { progress_interface(iter_counter, iter_conv, norm_iter_conv, loop_benchmark); }
92             break;
93     default:
94         break;
95     }
96 }
97
98 /**
99 * Prints the results of the K-means algorithm.
100 *
101 * @param[in] bench_results the benchmarking results.
102 * @param[in] iter_conv the final normalized K-means error.
103 *
104 * @see `Benchmark.h`
105 */
106 void kmeans_termination(const std::pair<std::chrono::duration<double>, std::chrono::duration<double>> bench_results, const long double iter_conv)
107 {
108     std::cout << std::endl
109     << std::setw(18) << "[TOTAL TIME] " << std::setw(15) << bench_results.first.count()
110     << " [IN SECONDS]" << std::endl
111     << std::setw(18) << "[K-MEANS BENCH] " << std::setw(15) << bench_results.second.count()
112     << " [IN SECONDS]" << std::endl
113     << std::setw(18) << "[CONVERGENCE] " << std::setw(15) << iter_conv << std::endl;
114 }

```

Κώδικας 3.2.14: Η κεφαλίδα Benchmark.h

```
1 /**
2  * Benchmark.h
3  *
4  * In this header file, we define some functions that help us
5  * benchmark some crucial parts of the K-means algorithm. Using
6  * these functions, we can benchmark the K-means runtime, as well as
7  * how long did it take for our implementation to converge.
8  */
9
10 #pragma once
11
12 #include "Common.h"
13
14 std::chrono::time_point<std::chrono::system_clock> init_benchmark(void);
15 std::chrono::time_point<std::chrono::system_clock> bench_convergence(void);
16 std::chrono::duration<double> bench_loop(std::chrono::time_point<std::chrono::system_clock>, std::chrono::time_point<std::chrono::system_clock>);
17 std::chrono::time_point<std::chrono::system_clock> terminate_bench(void);
18 std::pair<std::chrono::duration<double>, std::chrono::duration<double>> benchmark_results(std::chrono::time_point<std::chrono::system_clock> start, std::chrono::time_point<std::chrono::system_clock> loop_benchmark, std::chrono::time_point<std::chrono::system_clock> end);
```

Κώδικας 3.2.15: Το αρχείο Benchmark.cpp

```
1 #include "Benchmark.h"
2
3 /**
4  * Launches benchmark.
5  *
6  * @return the starting timepoint of the K-means algorithm.
7  */
8 std::chrono::time_point<std::chrono::system_clock> init_benchmark(void)
9 {
10     std::chrono::time_point<std::chrono::system_clock> start = std::chrono::system_clock::now();
11     return start;
12 }
13
14 /**
15  * Launches K-means optimization loop benchmark.
16  *
17  * @return the starting timepoint of the K-means optimization loop.
18  *
19  * @see kmeans.cpp
20  */
21 std::chrono::time_point<std::chrono::system_clock> bench_convergence(void)
22 {
23     std::chrono::time_point<std::chrono::system_clock> loop_benchmark = std::chrono::system_clock::now();
24     return loop_benchmark;
25 }
26
27 /**
28  * Computes per-loop benchmark.
29  *
30  * @param[in] startpoint loop starting time point
31  * @param[in] endpoint loop ending time point
32  *
33  * @return optimization loop benchmark.
34  *
35  * @see kmeans.cpp
36  */
37 std::chrono::duration<double> bench_loop(std::chrono::time_point<std::chrono::system_clock> startpoint, std::chrono::time_point<std::chrono::system_clock> endpoint)
38 {
39     return endpoint - startpoint;
40 }
41
42 /**
43  * Stops K-means benchmark.
44  *
45  * @return the ending timepoint of the K-means algorithm.
46  */
47 std::chrono::time_point<std::chrono::system_clock> terminate_bench(void)
48 {
49     std::chrono::time_point<std::chrono::system_clock> end = std::chrono::system_clock::now();
50     return end;
51 }
52
53 /**
54  * Computes K-means benchmarking results.
55  *
56  * @param[in] start the starting timepoint of the K-means algorithm.
57  * @param[in] loop_benchmark the starting timepoint of the K-means optimization loop.
58  * @param[in] end the ending timepoint of the K-means algorithm.
59  *
60  * @return a pair with the total execution time of the K-means algorithm in seconds
61  *         and the total time in seconds that the algorithm needed to converge.
62  */
63 std::pair<std::chrono::duration<double>, std::chrono::duration<double>> benchmark_results(
64     const std::chrono::time_point<std::chrono::system_clock> start,
65     const std::chrono::time_point<std::chrono::system_clock> loop_benchmark,
66     const std::chrono::time_point<std::chrono::system_clock> end)
67 {
68     std::chrono::duration<double> elapsed_seconds = end - start;
69     std::chrono::duration<double> k_means_benchmark = end - loop_benchmark;
70     std::pair<std::chrono::duration<double>, std::chrono::duration<double>> bench_res;
71     bench_res.first = elapsed_seconds;
72     bench_res.second = k_means_benchmark;
73     return bench_res;
74 }
```


Κώδικας 3.2.18: Η κεφαλίδα kmeans.h

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

```
/**
 * kmeans.h
 *
 * In this header file, we include all the
 * custom made header files used in the
 * K-means algorithm.
 */

#pragma once

#include "Vec.h"
#include "Center.h"
#include "Class.h"
#include "Optimize.h"
#include "Distance.h"
#include "Interface.h"
#include "Benchmark.h"
#include "Validation.h"
```

Κώδικας 3.2.19: Το αρχείο kmeans.cpp

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

```
#include "kmeans.h"

/**
 * Implements K-means clustering algorithm.
 *
 * @return 0, if the executable was terminated normally.
 */

int main(void)
{
    std::chrono::time_point<std::chrono::system_clock> loop_benchmark, prev_loop_benchmark, start;
    std::cout.precision(std::numeric_limits<float>::max_digits10);
    prev_loop_benchmark = start = init_benchmark();
    std::vector<std::array<float, Nv>> Vec;
    Vec.reserve(N);
    vec_init(Vec);
    multidimensional_float_vector_interface(Vec, "Vec");
    export_multidimensional_float_vector(Vec, "Vec");
    std::vector<std::array<float, Nv>> old_Center;
    old_Center.reserve(Nc);
    init_centers(Vec, old_Center);
    multidimensional_float_vector_interface(old_Center, "old_Center");
    export_multidimensional_float_vector(old_Center, "old_Center");
    std::vector<std::array<float, Nv>> new_Center;
    new_Center.reserve(Nc);
    long double iter_conv = std::numeric_limits<double>::infinity();
    long double norm_iter_conv = std::numeric_limits<double>::infinity();
    std::array<std::vector<int>, Nc> Classes;
    std::chrono::time_point<std::chrono::system_clock> optimization_benchmark = init_benchmark();
    for (int iter_counter = 1; norm_iter_conv > THR_KMEANS; iter_counter += 1)
    {
        std::chrono::duration<double> bench_res =
            bench_loop(prev_loop_benchmark, loop_benchmark = init_benchmark());
        kmeans_progress(iter_counter, iter_conv, norm_iter_conv, bench_res, VERBOSITY);
        track_kmeans_progress(iter_counter, norm_iter_conv);
        compute_classes(Vec, old_Center, Classes);
        multidimensional_int_array_interface(Classes, "Classes");
        export_multidimensional_integer_array(Classes, "Classes" + std::to_string(iter_counter));
        optimize_center(Vec, new_Center, Classes);
        multidimensional_float_vector_interface(new_Center, "new_Center");
        export_multidimensional_float_vector(new_Center, "new_Center" + std::to_string(iter_counter));
        long double prev_iter_conv = iter_conv;
        iter_conv = convergence(new_Center, old_Center);
        norm_iter_conv = normalize_convergence(iter_conv, prev_iter_conv, iter_counter);
        old_Center.swap(new_Center);
        multidimensional_float_vector_interface(old_Center, "old_Center");
    }
    std::chrono::time_point<std::chrono::system_clock> end = terminate_bench();
    std::pair<std::chrono::duration<double>, std::chrono::duration<double>> bench_res =
        benchmark_results(start, optimization_benchmark, end);
    kmeans_termination(bench_res, norm_iter_conv);
    export_kmeans_progress("kmeans_results");
    return 0;
}
```