



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение высшего
образования
«Московский государственный технический университет имени Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

Отчет по лабораторной работе №4 по курсу "Анализ алгоритмов"

Тема Многопоточные вычисления

Студент Романов С.К.

Группа ИУ7-55Б

Оценка (баллы)

Преподаватели Волкова Л.Л., Строганов Ю.В.

Москва — 2023 г.

Оглавление

ВВЕДЕНИЕ	2
1 Аналитическая часть	4
1.1 Алгоритм DBSCAN	4
1.2 Параллельный алгоритм DBSCAN	7
2 Конструкторская часть	10
2.1 Разработка алгоритмов	10
3 Технологическая часть	19
3.1 Требования к программному обеспечению	19
3.2 Средства реализации	19
3.3 Реализация алгоритмов	19
3.3.1 Последовательный алгоритм DBSCAN	20
3.3.2 Параллельный алгоритм DBSCAN	22
3.4 Тестовые данные	29
4 Исследовательская часть	30
4.1 Пример выполнения	30
4.2 Время выполнения алгоритмов	33
ЗАКЛЮЧЕНИЕ	36
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	37

ВВЕДЕНИЕ

Многопоточность — способность центрального процессора (CPU) или одного ядра в многоядерном процессоре одновременно выполнять несколько процессов или потоков, соответствующим образом поддерживаемых операционной системой.

Этот подход отличается от многопроцессорности, так как многопоточность процессов и потоков совместно использует ресурсы одного или нескольких ядер: вычислительных блоков, кэш-памяти ЦПУ или буфера перевода с преобразование.

Для распараллеливания может быть рассмотрена задача кластеризации пространственных данных.

Алгоритм DBSCAN (Density Based Spatial Clustering of Applications with Noise), был предложен Мартином Эстер, Хансом-Питером Кригель и коллегами в 1996 году как решение проблемы разбиения (изначально пространственных) данных на кластеры произвольной формы. Большинство алгоритмов, производящих плоское разбиение, создают кластеры по форме близкие к сферическим, так как минимизируют расстояние документов до центра кластера. Авторы DBSCAN экспериментально показали, что их алгоритм способен распознать кластеры различной формы.

Цель лабораторной работы:

- Описание и реализация параллельных вычислений

Задачи лабораторной работы:

1. Проанализировать последовательный и параллельный варианты алгоритма DBSCAN.
2. Определить средства программной реализации алгоритмов.
3. Реализовать разработанные алгоритмы.
4. Провести анализ затрат работы программы по времени, выяснить влияющие на них характеристики,
5. Создать отчёт, содержащий следующее:
 - актуальность исследования;
 - характеристики предложенной реализации алгоритмов (по времени);
 - результаты тестирования;
 - Выводы об эффективности алгоритмов, выведенных из их сравнения по времени.

В ходе работы будут затронуты следующие темы.

1. Параллелизация вычислений.
2. Алгоритм DBSCAN.
3. Многопоточность.

1 Аналитическая часть

1.1 Алгоритм DBSCAN

Идея, положенная в основу алгоритма, заключается в том, что внутри каждого кластера наблюдается типичная плотность точек (объектов), которая заметно выше, чем плотность снаружи кластера, а также плотность в областях с шумом ниже плотности любого из кластеров.[1]

Ещё точнее, что для каждой точки кластера её соседство заданного радиуса должно содержать не менее некоторого числа точек, это число точек задаётся пороговым значением. Перед изложением алгоритма дадим необходимые определения.

Определение 1. Eps-соседство точки p , обозначаемое как $N_\varepsilon(p)$, определяется как множество документов, находящихся от точки p на расстояния не более ε : $N_\varepsilon(p) = \{q \in D | dist(p, q) \leq \varepsilon\}$. Поиска точек, чья $N_\varepsilon(p)$ содержит хотя бы минимальное число точек (MinPt) не достаточно, так как точки бывают двух видов: ядровые и граничные.

Определение 2. Точка p непосредственно плотно-достижима из точки q (при заданных ε и MinPt), если $p \in N_\varepsilon(q)$ и $|N_\varepsilon(q)| \geq MinPt$ (рис. 1.1)

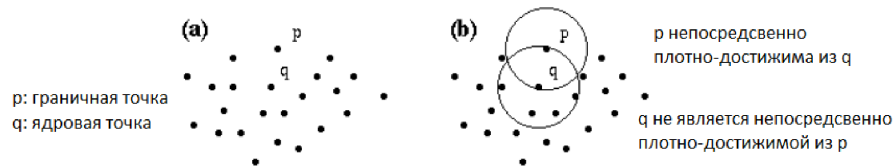


Рисунок 1.1 – Пример точек, находящихся в отношении непосредственно плотной достижимости

Определение 3. Точка p плотно-достижима из точки q (при заданных ε и MinPt), если существует последовательность точек $q = p_1, p_2, \dots, p_n = p : p_{i+1}$ непосредственно плотно-достижимы из p_i . Это отношение транзитивно, но не симметрично в общем случае, однако симметрично для двух ядровых точек.

Определение 4. Точка p плотно-связана с точкой q (при заданных ε и MinPt), если существует точка o : p и q плотно-достижимы из o (при заданных ε и MinPt), (рис. 1.2)

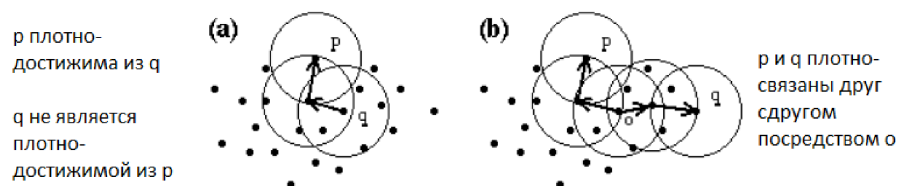


Рисунок 1.2 – Пример точек, находящихся в отношении плотной связанности.

Теперь мы готовы дать определения кластеру и шуму.

Определение 5. Кластер C_j (при заданных ε и MinPt) – это не пустое подмножество документов, удовлетворяющее следующим условиям:

1. $\forall p, q$: если $p \in C_j$ и q плотно-достижима из p (при заданных ε и MinPt), то $q \in C_j$,
2. $\forall p, q \in C_j$: p плотно связана с q (при заданных ε и MinPt);

Итак, кластер – это множество плотно-связанных точек. В каждом кластере содержится хотя бы MinPt документов.

Шум – это подмножество документов, которые не принадлежат ни одному кластеру: $\{p \in D | \forall j : p \notin C_j, j = \overline{1, |\mathcal{C}|}\}$

Алгоритм DBSCAN для заданных значений параметров ε и MinPt исследует кластер следующим образом: сначала выбирает случайную точку, являющуюся ядровой, в качестве затравки, затем помещает в кластер саму затравку и все точки, плотно-достижимые из неё.

Алгоритм в общем виде.

DBSCAN

Вход: множество точек документов \mathcal{D} , ε и MinPt .

Выход: множество кластеров $\mathcal{C} = \{C_j\}$.

Шаг 1. Установить всем элементам множества \mathcal{D} флаг «не посещён».

Присвоить текущему кластеру C_j нулевой номер, $j := 0$.

Множество шумовых документов $\text{Noise} := \emptyset$.

Шаг 2. Для каждого $d_i \in \mathcal{D}$ такого, что $\text{флаг}(d_i) = \text{«не посещён»}$, выполнить:

Шаг 3. $\text{флаг}(d_i) := \text{«посещён»}$;

Шаг 4. $N_i := N_\varepsilon(d_i) = \{q \in \mathcal{D} \mid \text{dist}(d_i, q) \leq \varepsilon\}$

Шаг 5. Если $|N_i| < \text{MinPt}$, то
 $\text{Noise} := \text{Noise} + \{d_i\}$

 иначе

 номер следующего кластера $j := j + 1$;

$\text{EXPANDCLUSTER}(d_i, N_i, \mathcal{C}, \varepsilon, \text{MinPt})$;

EXPANDCLUSTER

Вход: текущая точка d_i , его ε -соседство N_i , текущий кластер C_j и ε , MinPt .

Выход: кластер C_j

Шаг 1. $C_j := C_j + \{d_i\}$;

Шаг 2. Для всех документов $d_k \in N_i$:

Шаг 3. Если $\text{флаг}(d_k) = \text{«не посещён»}$, то

Шаг 4. $\text{флаг}(d_k) := \text{«посещён»}$;

Шаг 5. $N_{ik} := N_\varepsilon(d_k)$;

Шаг 6. $|N_{ik}| \geq \text{MinPt}$, то $N_i := N_i + N_{ik}$;

Шаг 7. Если $\nexists p : d_k \in C_p, p = \overline{1, |\mathcal{C}|}$, то $C_j := C_j + \{d_k\}$;

1.2 Параллельный алгоритм DBSCAN

Параллельный алгоритм DBSCAN заключается в том, что точки помещаются в непересекающиеся d -мерные клетки с длиной стороны ε/\sqrt{d} по их координатам (рис. 1.3(b)). Ячейки обладают тем свойством, что все точки внутри клетки находятся на расстоянии ε друг от друга и будут принадлежать к одному и тому же кластеру.[2] Затем (рис. 1.3(c)) мы отмечаем основные точки, после чего мы генерируем кластеры для основных точек следующим образом:

Мы создаем граф содержащий по одной вершине на ядровую ячейку (ячейка, содержащая не менее одной ядровой точки), и соединить две вершины, если ближайшая пара ядровых точек из двух ячеек находится на расстоянии ε . Далее будем ссылаться на него как на клеточный граф. Сам этот шаг проиллюстрирован на рисунке 1.3(d). Затем находим компоненты связности графа ячеек для назначения идентификаторов кластеров точкам в ядровых ячейках и назначаем идентификаторы кластеров для граничных точек. В конце возвращаем список кластеров.

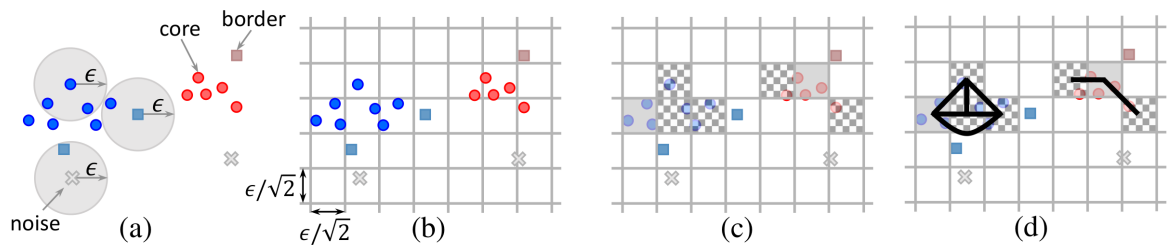


Рисунок 1.3 – Параллельный алгоритм DBSCAN

Алгоритм в общем виде.

DBSCAN

Вход: множество индексированных точек \mathcal{P} , ε и MinPt .

Выход: множество кластеров $\mathcal{C} = \{C_j\}$.

Шаг 1. Разбить множество \mathcal{P} точек на множество ячеек \mathcal{G}

Шаг 2. $\text{coreFlags} := \text{MarkCore}(\mathcal{G}, \mathcal{P}, \varepsilon, \text{MinPt})$

Шаг 3. $\mathcal{C} := \text{ClusterCore}(\mathcal{G}, \mathcal{P}, \text{coreFlags}, \varepsilon, \text{MinPt})$

Шаг 4. $\text{ClusterBorder}(\mathcal{G}, \mathcal{P}, \text{coreFlags}, \mathcal{C}, \varepsilon, \text{MinPt})$

Шаг 5. Вернуть \mathcal{C}

MarkCore

Вход: множество индексированных точек \mathcal{P} , множество ячеек \mathcal{G} , ε и MinPt .

Выход: список ядровых точек coreFlags .

1. $\text{coreFlags} := \{0, \dots, 0\}$ $\backslash \backslash$ Длина равна размеру \mathcal{P}
 2. Параллельно $\forall g \in \mathcal{G}$
 3. Если $\text{len}(g) \geq \text{minPt}$
 4. Для $\forall p \in g$
 5. $\text{coreFlags}[p] := 1$
 6. иначе
 7. для $\forall p \in g$
 8. $\text{count} := \text{len}(g)$
 9. для $\forall h \in g.\text{NeighborCells}(\varepsilon)$
 10. $\text{count} := \text{count} + \text{RangeCount}(p, \varepsilon, h)$
 11. Если $\text{count} \geq \text{minPt}$
 12. $\text{coreFlags}[p] := 1$
 13. Вернуть coreFlags
-

ClusterCore

Вход: множество индексированных точек \mathcal{P} , множество ячеек \mathcal{G} ,
флаги ядровых ячеек $coreFlags$, ε и $MinPt$.

Выход: Кластеры \mathcal{C} .

1. $uf := UnionFind()$ $\backslash\backslash$ Инициализация Union-find структуры
 2. Параллельно $\forall g \in \mathcal{G}$
 3. Для $\forall h \in g.NeighborCells()$
 4. Если $g > h$ $\&\&$ $uf.Find(g) \neq uf.Find(h)$
 5. Если $Connected(g, h)$
 6. $uf.Link(g, h)$
 7. для $\forall p \in g$
 8. $\mathcal{C} := \{-1, \dots, -1\}$ Длина множества \mathcal{P}
 9. Параллельно $\forall g \in \mathcal{G}$
 10. $\forall p \in g : coreFlags[p] = 1$
 11. $\mathcal{C}[p] := uf.Find(g)$
 12. Вернуть \mathcal{C}
-

ClusterBorder

Вход: множество индексированных точек \mathcal{P} , множество ячеек \mathcal{G} ,
кластеры \mathcal{C} , флаги ядровых ячеек $coreFlags$, ε и $MinPt$.

Выход: Обновленные кластеры \mathcal{C}

1. Параллельно $\forall g \in \mathcal{G} : len(g) < minPt$
 2. Для $\forall p \in cell\ g : coreFlags[p] = 0$
 3. Для $\forall h \in g \cup g.NeighborCells(\varepsilon)$
 4. Для $\forall q \in cell\ h : coreFlags[q] = 1$
 5. Если $dist(p, q) \leq \varepsilon$
 6. $\mathcal{C}[p] := \mathcal{C}[p] \cup \mathcal{C}[q]$
 7. Вернуть \mathcal{C}
-

Вывод

В данном разделе были рассмотрены алгоритмы DBSCAN и его распараллеленного аналога.

2 Конструкторская часть

2.1 Разработка алгоритмов

На рисунках 2.1, 2.2 представлена схема линейного алгоритма DBSCAN.

На рисунках 2.3, 2.4, 2.5, 2.6, 2.7, 2.8 представлена схема параллельного алгоритма DBSCAN.

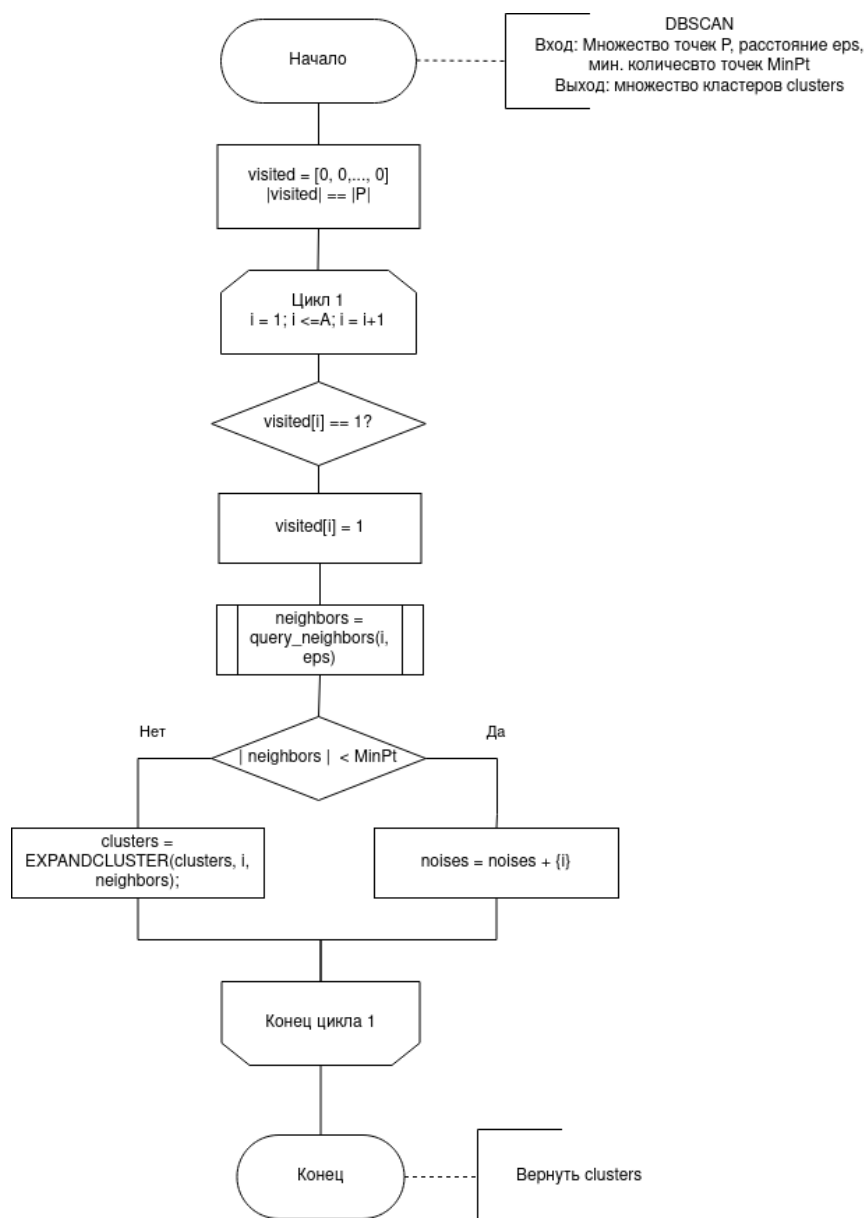


Рисунок 2.1 – Линейный алгоритм DBSCAN

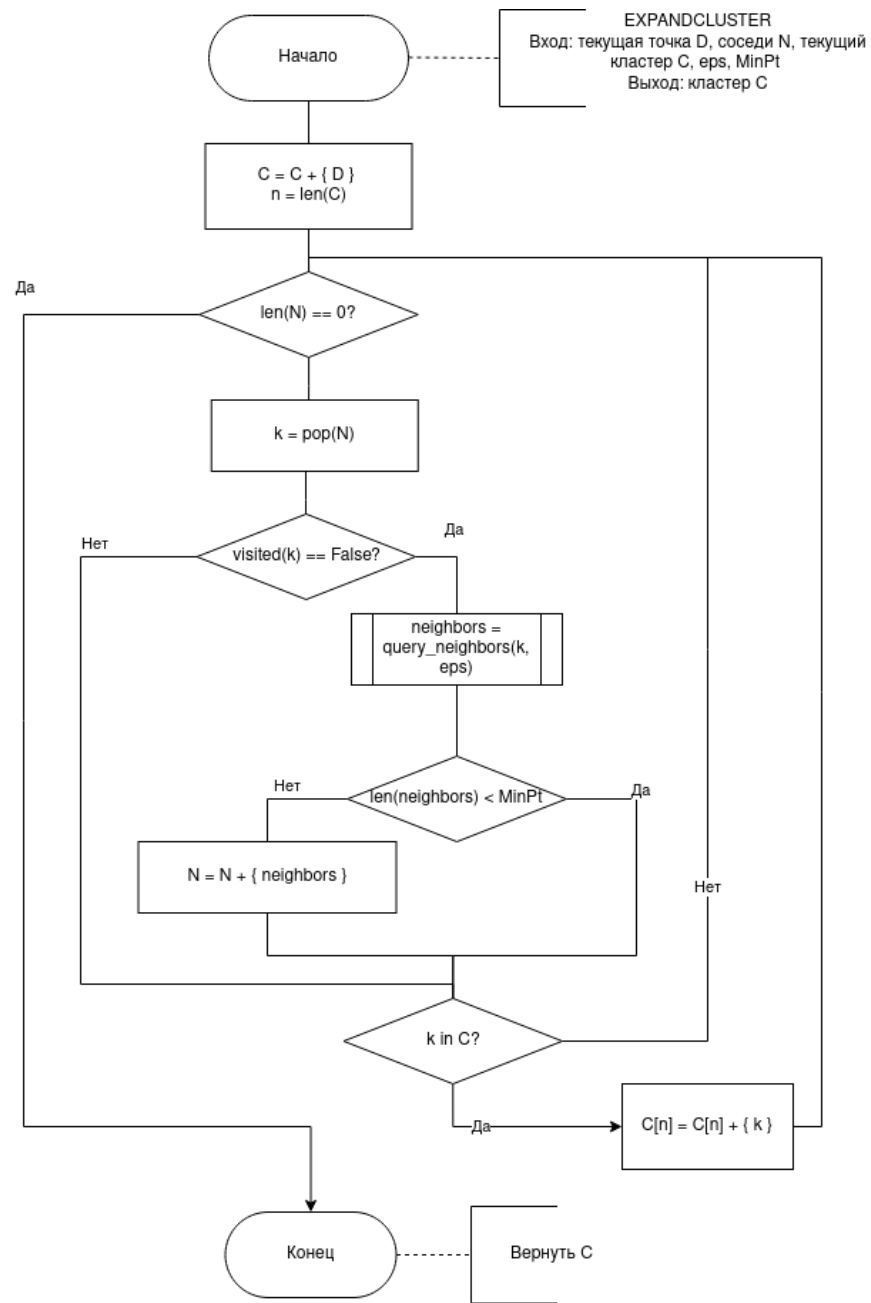


Рисунок 2.2 – Функция EXPANDCLUSTERS линейного алгоритма DBSCAN

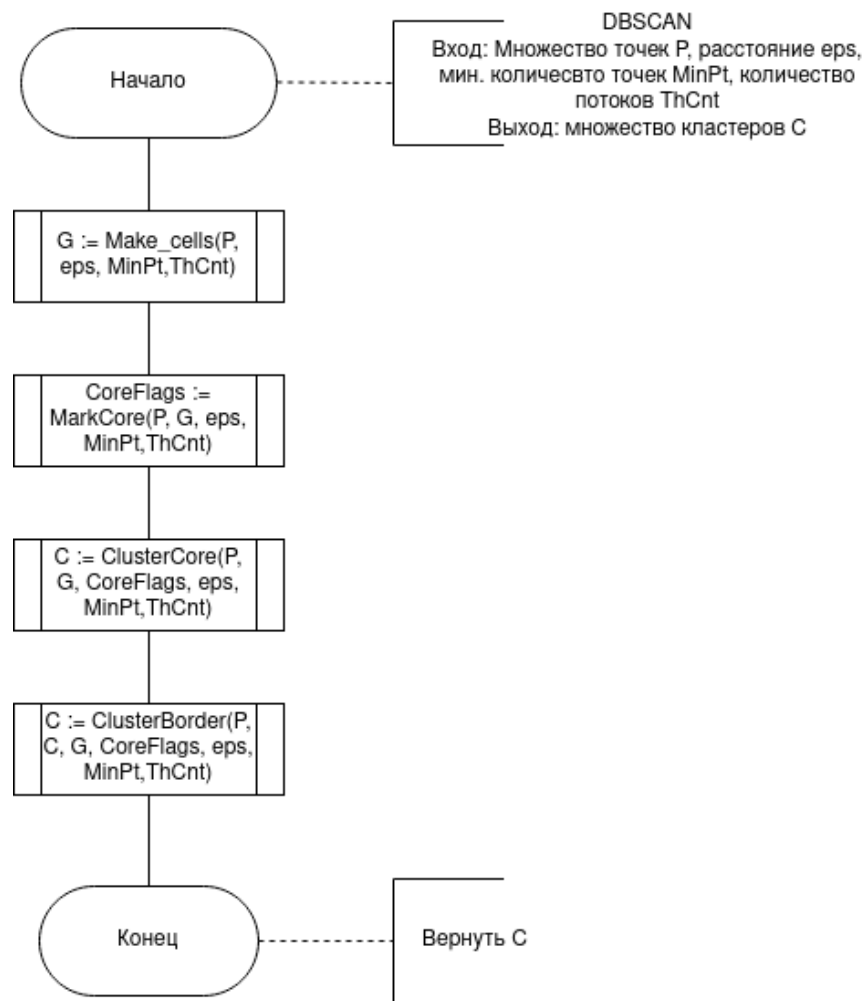


Рисунок 2.3 – Параллельный алгоритм DBSCAN

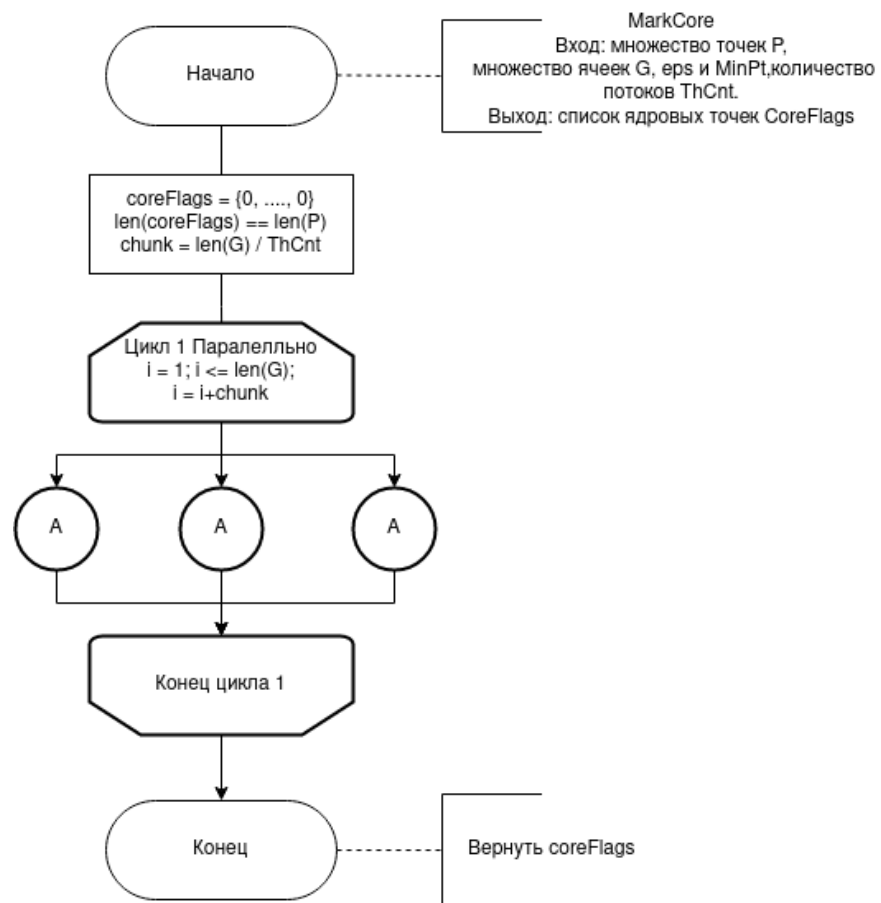


Рисунок 2.4 – Функция markCore параллельного алгоритма DBSCAN, часть 1

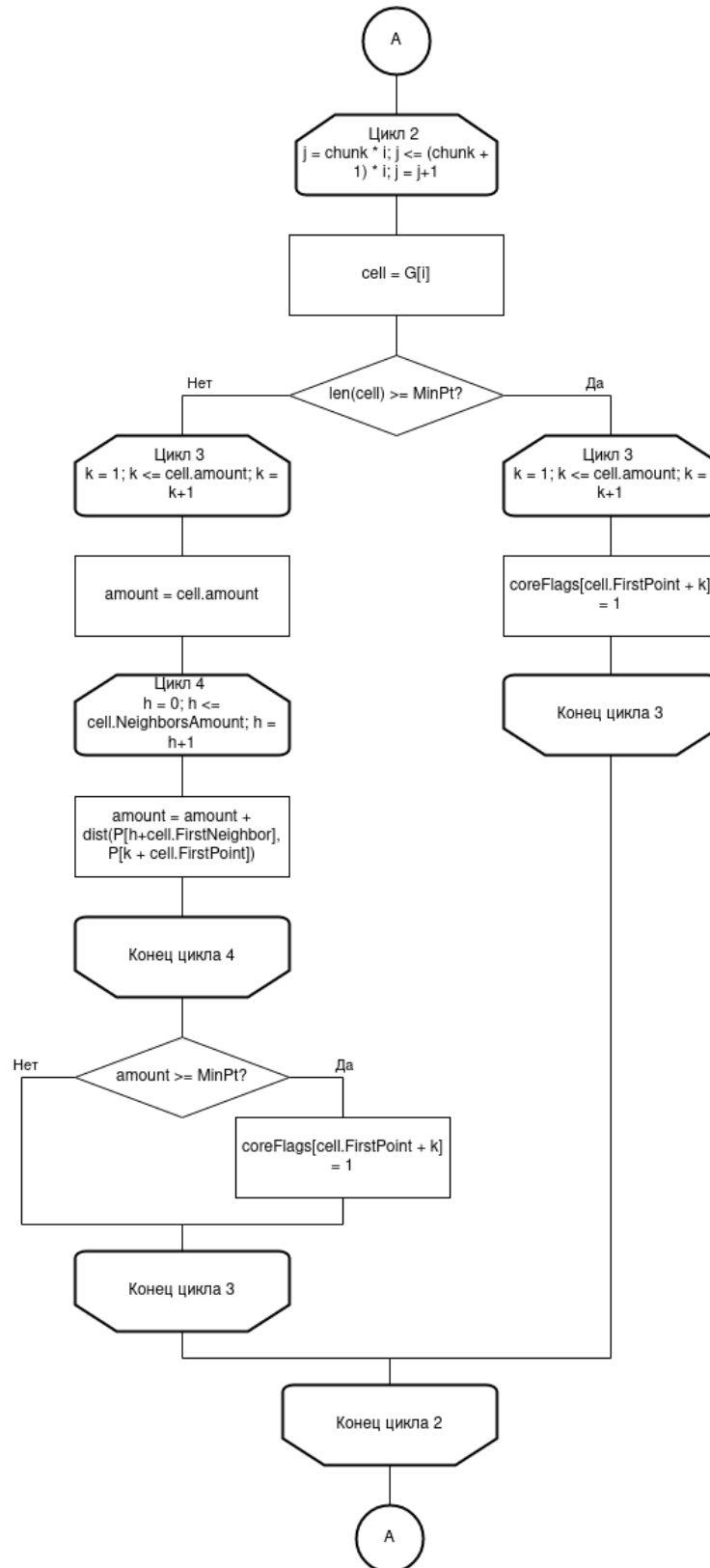


Рисунок 2.5 – Функция markCore параллельного алгоритма DBSCAN, часть 2

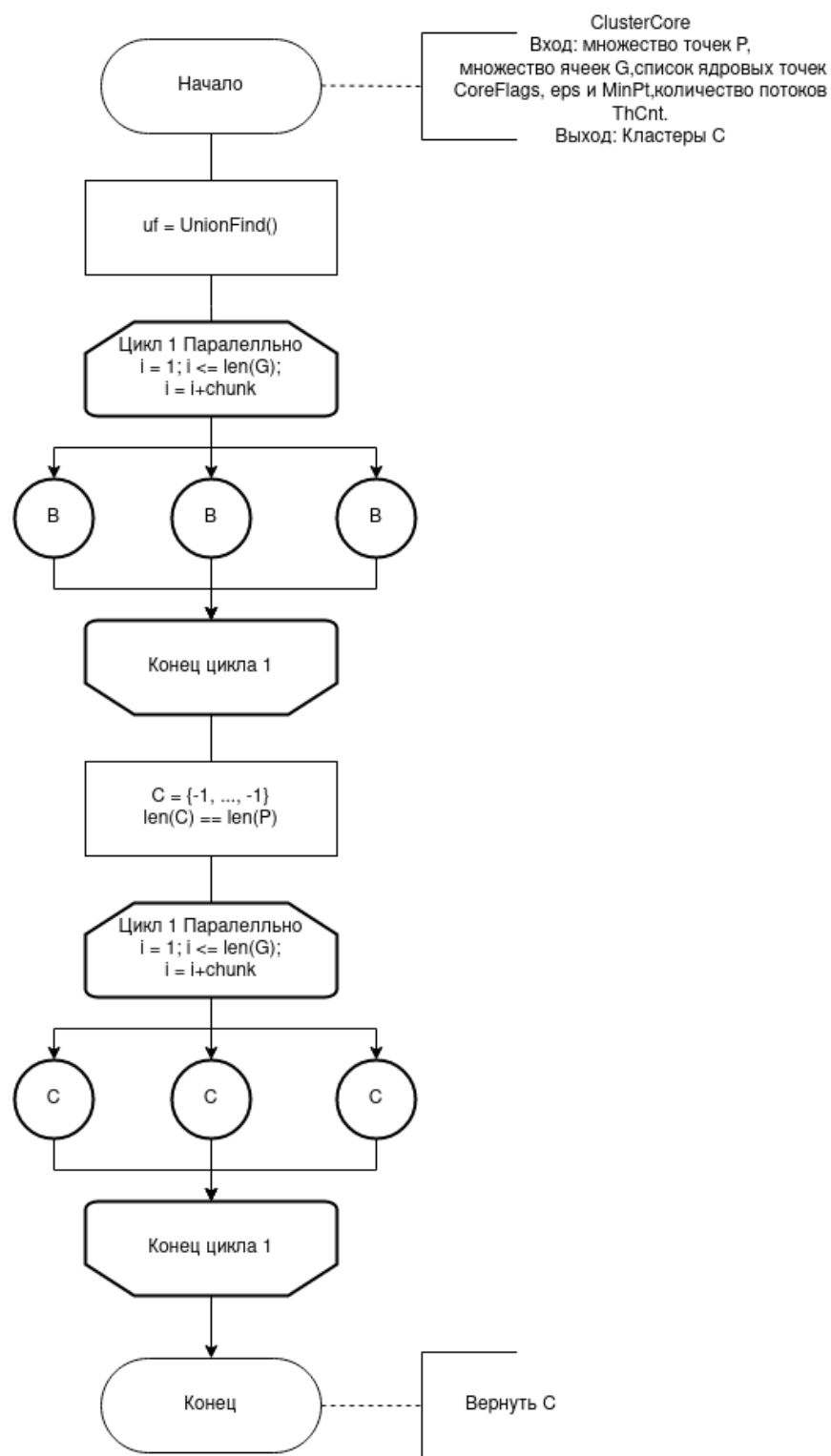


Рисунок 2.6 – Функция clusterCore параллельного алгоритма DBSCAN, часть 1

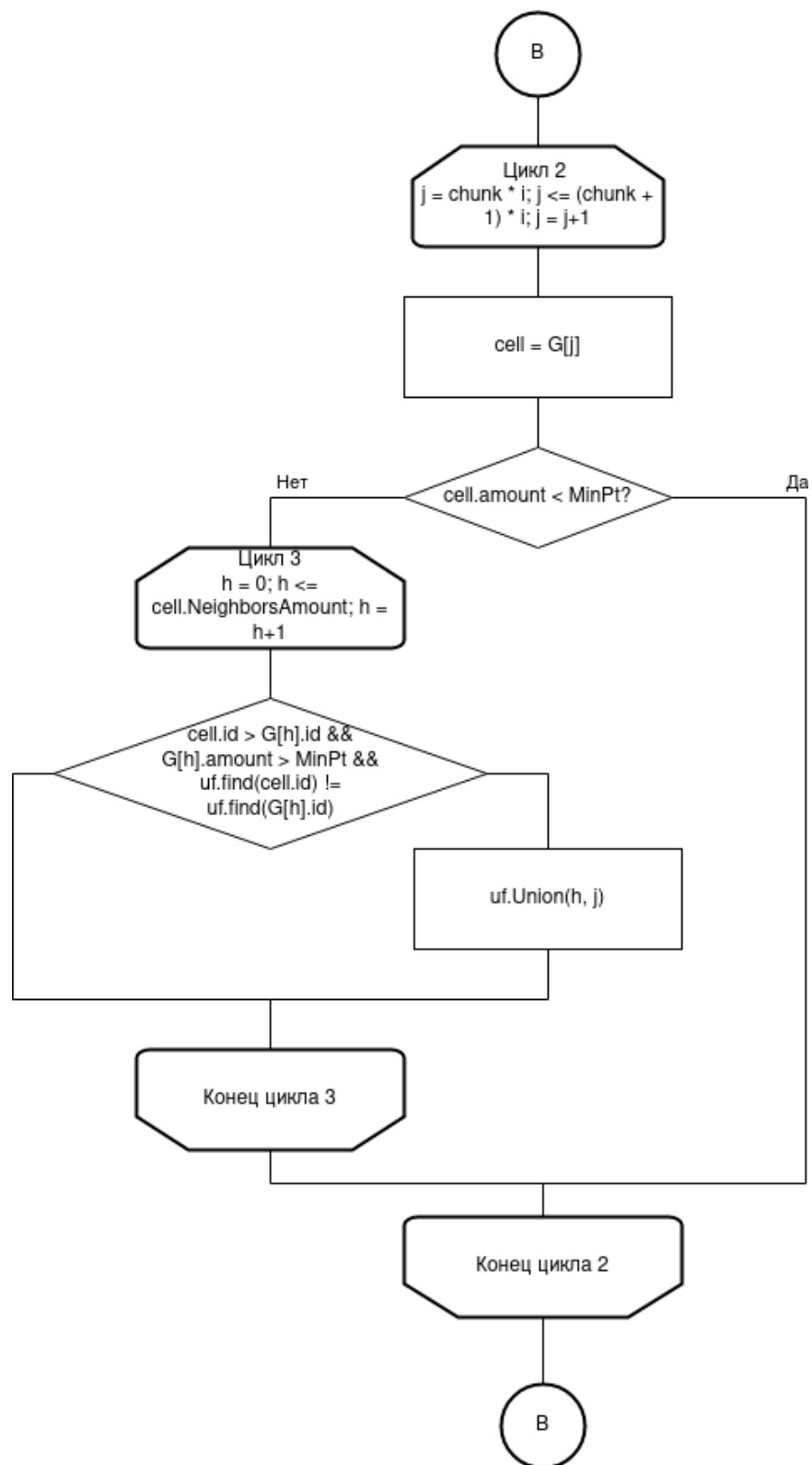


Рисунок 2.7 – Функция clusterCore параллельного алгоритма DBSCAN, часть 2

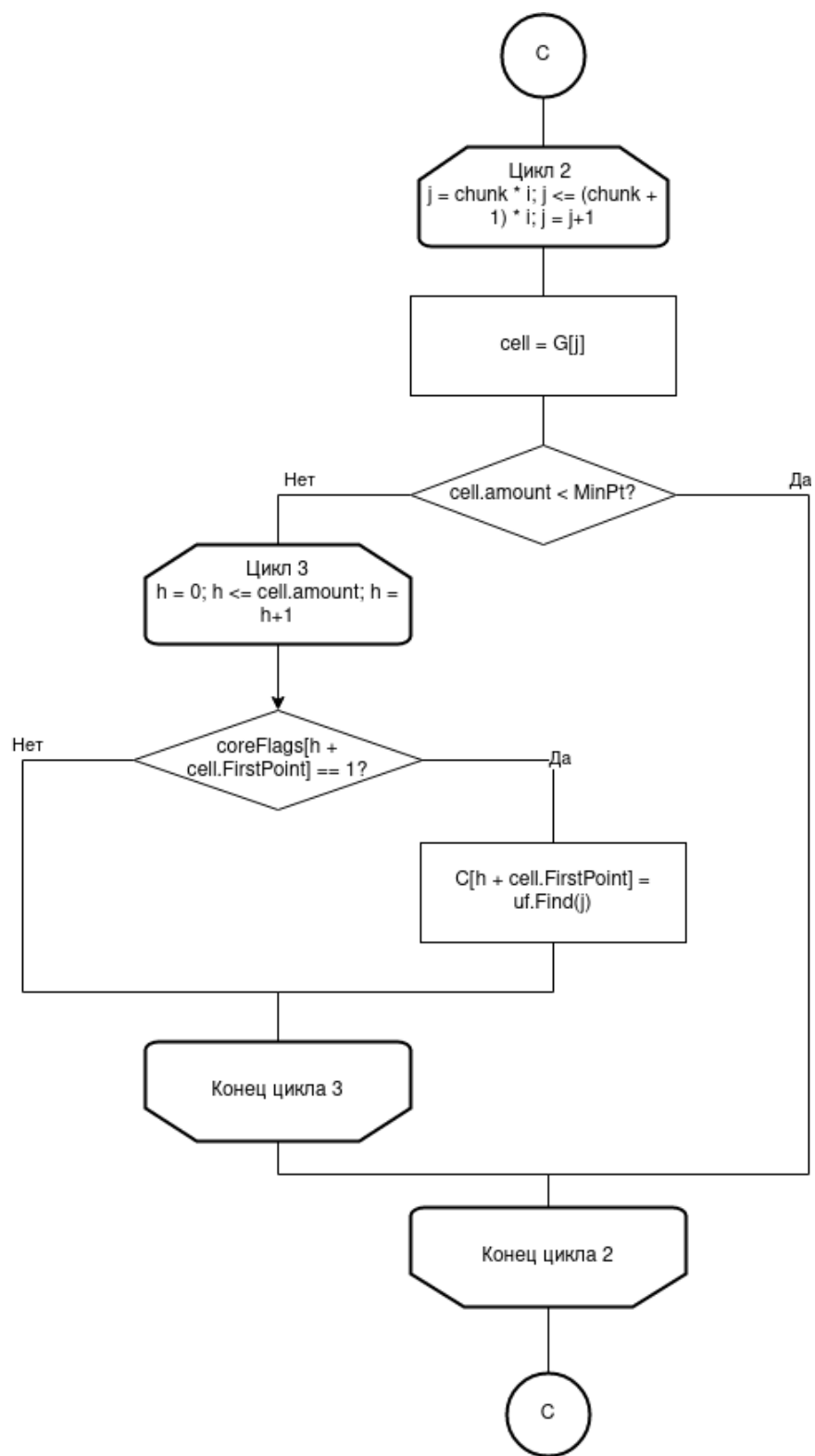


Рисунок 2.8 – Функция clusterCore параллельного алгоритма DBSCAN, часть 3

Вывод

На основе теоретических данных, полученных из аналитического раздела, были построены схемы алгоритмов DBSCAN, последовательного и параллельного.

3 Технологическая часть

В данном разделе приведены требования к программному обеспечению, средства реализации и сами реализации алгоритмов.

3.1 Требования к программному обеспечению

К программе предъявляется ряд условий:

- на вход подается путь к файлу, содержащий некие двумерные данные, а также числа ϵ и MinPt , определяющие характеристики алгоритма;
- на выход ПО должно выводить результат алгоритма DBSCAN в графическом виде;

3.2 Средства реализации

Для реализации данной лабораторной работы необходимо установить следующее программное обеспечение:

- Rust Programming Language v1.64.0 — язык программирования [3]
- Criterion.rs v0.4.0 — средство визуализации данных
- LaTeX — система верстки документов

3.3 Реализация алгоритмов

В следующих листингах представлены следующие алгоритмы:

1. В листингах 3.1 и 3.2 представлен последовательный алгоритм DBSCAN.
2. В листингах 3.3, 3.4, 3.5, 3.6 и 3.7 представлен параллельный алгоритм DBSCAN.

3.3.1 Последовательный алгоритм DBSCAN

Листинг 3.1 – Последовательный алгоритм DBSCAN

```
1  pub fn run(&mut self)
2  {
3      let len = self.dataset.lock().unwrap().labels_amount();
4      for col_a in 0..len {
5          for col_b in col_a + 1..len {
6              self.focus = Focus(col_a, col_b);
7              self.visited = vec![false; self.dataset.lock().unwrap().records_amount()];
8              self.clusters.push((vec![], col_a, col_b));
9              for i in 0..self.visited.len() {
10                 if self.visited[i] {
11                     continue;
12                 }
13                 self.visited[i] = true;
14                 let mut neighbors = self.dataset.lock().unwrap().query_neighbors(
15                     i,
16                     self.eps,
17                     self.focus.clone().into(),
18                 );
19
20                 if neighbors.len() < self.min_pts {
21                     self.noises.push(i);
22                 } else {
23                     self.expand_cluster(i, &mut neighbors);
24                 }
25             }
26         }
27     }
28 }
```

Листинг 3.2 – Последовательный алгоритм DBSCAN, расширение кластеров

```
1  fn expand_cluster(&mut self, index: usize, neighbors: &mut Vec<usize>)
2  {
3      let n = self.clusters.len() - 1;
4      self.clusters[n].0.push(vec![index]);
5      let len = self.clusters[n].0.len() - 1;
6      while !neighbors.is_empty() {
7          let k = neighbors.pop().unwrap();
8          if !self.visited[k] {
9              self.visited[k] = true;
10             let new_neighbors = self.dataset.lock().unwrap().query_neighbors(
11                 k,
12                 self.eps,
13                 self.focus.clone().into(),
14             );
15             if new_neighbors.len() >= self.min_pts {
16                 neighbors.extend(new_neighbors);
17             }
18         }
19         if self.not_in_clusters(k) {
20             self.clusters[n].0[len].push(k);
21         }
22     }
23 }
24
25 fn not_in_clusters(&self, index: usize) -> bool
26 {
27     let mut res = true;
28     let n = self.clusters.len() - 1;
29     for cluster in self.clusters[n].0.iter() {
30         if cluster.contains(&index) {
31             res = false;
32             break;
33         }
34     }
35
36     res
37 }
```

3.3.2 Параллельный алгоритм DBSCAN

Листинг 3.3 – Параллельный алгоритм DBSCAN

```
1  pub fn run(&mut self)
2  {
3      let len = self.dataset.lock().unwrap().labels_amount();
4      for col_a in 0..len {
5          for col_b in col_a + 1..len {
6              self.focus = Focus(col_a, col_b);
7              self.make_cells();
8              #[cfg(debug_assertions)]
9              {
10                 draw_cells(self);
11             }
12             self.mark_core();
13             self.cluster_core();
14             self.cluster_border();
15         }
16     }
17 }
```

Листинг 3.4 – Параллельный алгоритм DBSCAN, получение ячеек

```
1  fn make_cells(&mut self)
2  {
3      let binding = self.dataset.lock().unwrap();
4
5      let vx = &binding.get_data()[self.focus.0];
6      let vy = &binding.get_data()[self.focus.1];
7      self.pairs = vx
8          .iter()
9          .cloned()
10         .zip(vy.iter().cloned())
11         .collect::<Vec<_>>();
12
13     parallel_sort(&mut self.pairs, &pairs_compare_by_x, self.threads_cnt);
14
15     let minx = self.pairs[0].0;
16     let mut cur_strip = 0;
17     self.strips.push(Strip {
18         id: 0,
19         amount: 0,
20         point_id: 0,
21         cells_amount: 0,
22         cell_id: 0,
23     });
24     let eps_len = self.eps / 2f64.sqrt();
25     for i in 0..self.pairs.len() {
26         if self.pairs[i].0 > (cur_strip + 1) as f64 * (eps_len) + minx {
27             parallel_sort(
28                 &mut self.pairs[self.strips[cur_strip].point_id..i],
29                 &pairs_compare_by_y,
30                 self.threads_cnt,
31             );
32             cur_strip += 1;
33             self.strips.push(Strip {
34                 id: cur_strip,
35                 amount: 0,
36                 point_id: i,
37                 cells_amount: 0,
38                 cell_id: 0,
39             });
40         }
41         self.strips[cur_strip].amount += 1;
42     }
43
44     let mut cell_id = 0;
```



```

45     for strip in &mut self.strips {
46         self.cells.push(Cell {
47             id: cell_id,
48             amount: 0,
49             point_id: strip.point_id,
50             strip_id: strip.id,
51             neighbors: vec![],
52         });
53         strip.cell_id = cell_id;
54         strip.cells_amount = 1;
55         for i in strip.point_id..(strip.point_id + strip.amount) {
56             if self.pairs[i].1 > self.pairs[self.cells[cell_id].point_id].1 + eps_len {
57                 cell_id += 1;
58                 self.cells.push(Cell {
59                     id: cell_id,
60                     amount: 0,
61                     point_id: i,
62                     strip_id: strip.id,
63                     neighbors: vec![],
64                 });
65                 strip.cells_amount += 1;
66             }
67             self.cells[cell_id].amount += 1;
68         }
69         cell_id += 1;
70     }
71     drop(binding);
72     self.make_cell_neighbors();
73 }

```

Листинг 3.5 – Параллельный алгоритм DBSCAN, получение coreFlags

```

1  fn mark_core(&mut self) {
2      let mut visited = vec![false; self.pairs.len()];
3      let n = self.cells.len();
4      let chunk_size = n / self.threads_cnt;
5      let self_arc = Arc::new(&self);
6      let _ = crossbeam::scope(|scope| {
7          let mut visited = visited.iter_mut();
8          for slice in (0..n).step_by(chunk_size) {
9              let self_arc = self_arc.clone();
10             let mut visited = visited.by_ref()
11                 .take(
12                     -(self.cells[slice].point_id as i64)
13                     + (self.cells[std::cmp::min(slice + chunk_size, n) - 1].point_id
14                        as i64)
15                     + (self.cells[std::cmp::min(slice + chunk_size, n) - 1].amount as i64))
16                     as usize,
17                 )
18                 .collect::<Vec<_>>();
19             let k = self.cells[slice].point_id;
20             scope.spawn(move |_| {
21                 for i in slice..std::cmp::min(slice + chunk_size, n) {
22                     let cell = &self_arc.cells[i];
23                     if cell.amount >= self_arc.min_pts {
24                         for i in cell.point_id..(cell.point_id + cell.amount - 1) {
25                             *visited[i - k] = true;
26                         }
27                     } else {
28                         for i in cell.point_id..(cell.point_id + cell.amount - 1) {
29                             let mut amount = cell.amount;
30                             for neighbor in &cell.neighbors {
31                                 amount += self_arc.range_count(self_arc.pairs[i], *neighbor);
32                             }
33                             if amount >= self_arc.min_pts {
34                                 *visited[i - k] = true;
35                             }
36                         }
37                     }
38                 }
39             });
40         }
41     });
42     self.visited = visited.clone();
43 }

```

Листинг 3.6 – Параллельный алгоритм DBSCAN, получение кластеров

```

1  fn cluster_core(&mut self)
2  {
3      let mut union_find = UnionFind::new(self.cells.len());
4      let mut clusters: Vec<i32> = vec![-1; self.pairs.len()];
5      let n = self.cells.len();
6      let chunk_size = n / self.threads_cnt;
7      let self_arc = Arc::new(&self);
8      let union_find = Arc::new(Mutex::new(union_find));
9      let _ = crossbeam::scope(|scope| {
10         for slice in (0..n).step_by(chunk_size) {
11             let self_arc = self_arc.clone();
12             let union_find = union_find.clone();
13             scope.spawn(move |_| {
14                 for i in slice..std::cmp::min(slice + chunk_size, n) {
15                     let cell = &self_arc.cells[i];
16                     if cell.amount >= self_arc.min_pts {
17                         for neighbor in &cell.neighbors {
18                             let first = union_find.lock().unwrap().find(cell.id);
19                             let second = union_find
20                                 .lock()
21                                 .unwrap()
22                                 .find(self_arc.cells[*neighbor].id);
23                             if cell.id > self_arc.cells[*neighbor].id
24                                 && self_arc.cells[*neighbor].amount >= self_arc.min_pts
25                                 && first != second
26                             {
27                                 union_find.lock().unwrap().union(i, *neighbor);
28                             }
29                         }
30                     }
31                 }
32             });
33         }
34     });
35     #[cfg(debug_assertions)]
36     {
37         println!("union find done");
38     }
39     let _ = crossbeam::scope(|scope| {
40         let mut clusters = clusters.iter_mut();
41         for (i, _) in (0..n).step_by(chunk_size).enumerate() {
42             let self_arc = self_arc.clone();
43             let union_find = union_find.clone();
44             let mut clusters = clusters

```

```

45         .by_ref()
46         .take(
47             (-(self.cells[i * chunk_size].point_id as i64)
48                 + (self.cells[std::cmp::min((i + 1) * chunk_size, n) - 1].point_id
49                     as i64)
50                 + (self.cells[std::cmp::min((i + 1) * chunk_size, n) - 1].amount
51                     as i64)) as usize,
52         )
53         .collect::<Vec<_>>());
54     let k = self.cells[i * chunk_size].point_id;
55     scope.spawn(move |_| {
56         for j in i * chunk_size..std::cmp::min((i + 1) * chunk_size, n) {
57             let cell = &self_arc.cells[j];
58             if cell.amount >= self_arc.min_pts {
59                 for p in cell.point_id..(cell.point_id + cell.amount) {
60                     if self_arc.visited[p] {
61                         *clusters[p - k] = union_find.lock().unwrap().find(j) as i32;
62                     }
63                 }
64             }
65         }
66     });
67 }
68 });
69 #[cfg(debug_assertions)]
70 {
71     println!("clusters done");
72 }
73 self.clusters = clusters.to_vec();
74 }

```

Листинг 3.7 – Параллельный алгоритм DBSCAN, расширение кластеров

```

1      let n = self.cells.len();
2      let chunk_size = n / self.threads_cnt;
3      let self_arc = Arc::new(&self);
4      let self_mutex = Arc::new(Mutex::new(self.clusters.clone()));
5      let _ = crossbeam::scope(|scope| {
6          for slice in (0..n).step_by(chunk_size) {
7              let self_arc = self_arc.clone();
8              let self_mutex = self_mutex.clone();
9              scope.spawn(move |_| {
10                 let cell = &self_arc.cells[slice];
11                 if cell.amount < self_arc.min_pts {
12                     for i in cell.point_id..(cell.point_id + cell.amount) {
13                         for neighbor in &cell.neighbors {
14                             for p in self_arc.cells[*neighbor].point_id
15                                 ..(self_arc.cells[*neighbor].point_id
16                                    + self_arc.cells[*neighbor].amount)
17                             {
18                                 if self_arc.clusters[p] != -1
19                                 && self_arc.dst(self_arc.pairs[i], self_arc.pairs[p])
20                                    <= self_arc.eps
21                                 {
22                                     self_mutex.lock().unwrap()[i] = self_arc.clusters[p];
23                                     break;
24                                 }
25                             }
26                         }
27                     }
28                 }
29             });
30         }
31     });
32     self.clusters = self_mutex.lock().unwrap().clone();
33 }
34
35 fn dst(&self, (x1, y1): (f64, f64), (x2, y2): (f64, f64)) -> f64

```

3.4 Тестовые данные

Одним из недостатков алгоритма DBSCAN является его неоднозначность, т.е. нет гарантии в однозначности результата. В связи с этим провести сравнение с неким эталонным результатом не представляется возможным

Вывод

В данном разделе был продемонстрирован исходный код алгоритма DBSCAN, последовательного и параллельного.

4 Исследовательская часть

Ниже приведены технические характеристики устройства, на котором было проведено тестирование ПО.

- Операционная система: Arch Linux [4] 64-bit.
- Оперативная память: 16 Гбайт, DDR4.
- Процессор: 11th Gen Intel® Core™ i5-11320H @ 3.20 ГГц [5]. 4 - физических ядра, 8 - логических.

4.1 Пример выполнения

На рис. 4.1 приведен пример выполнения реализации последовательного алгоритма, на рис. 4.2 – 4.4 — параллельного.

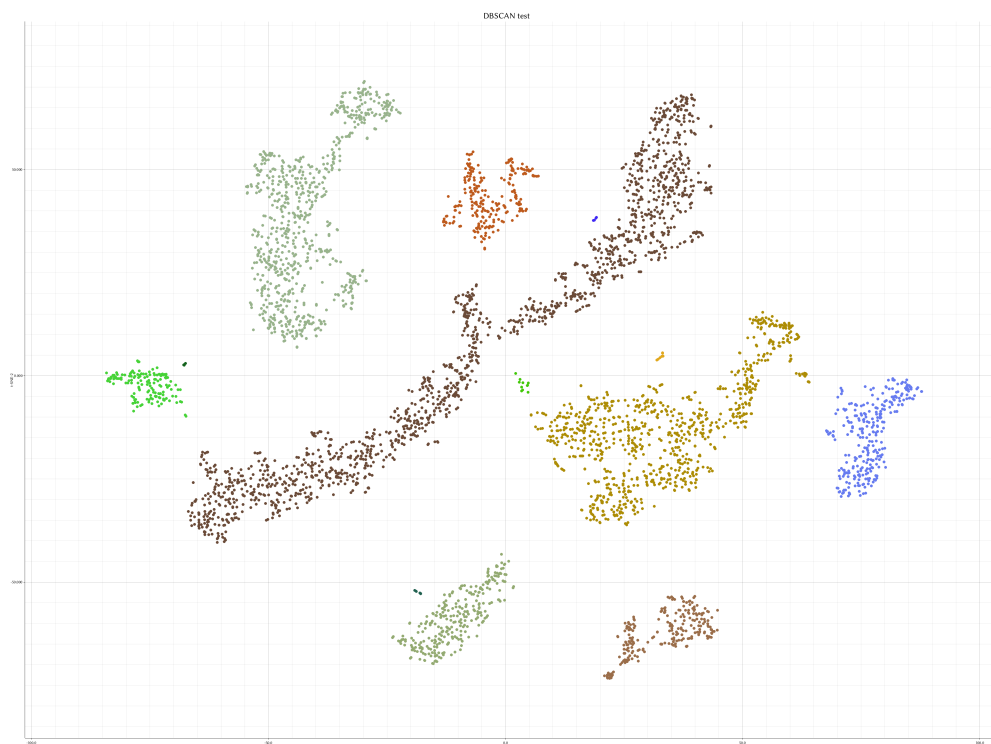


Рисунок 4.1 – Пример выполнения DBSCAN, последовательный алгоритм

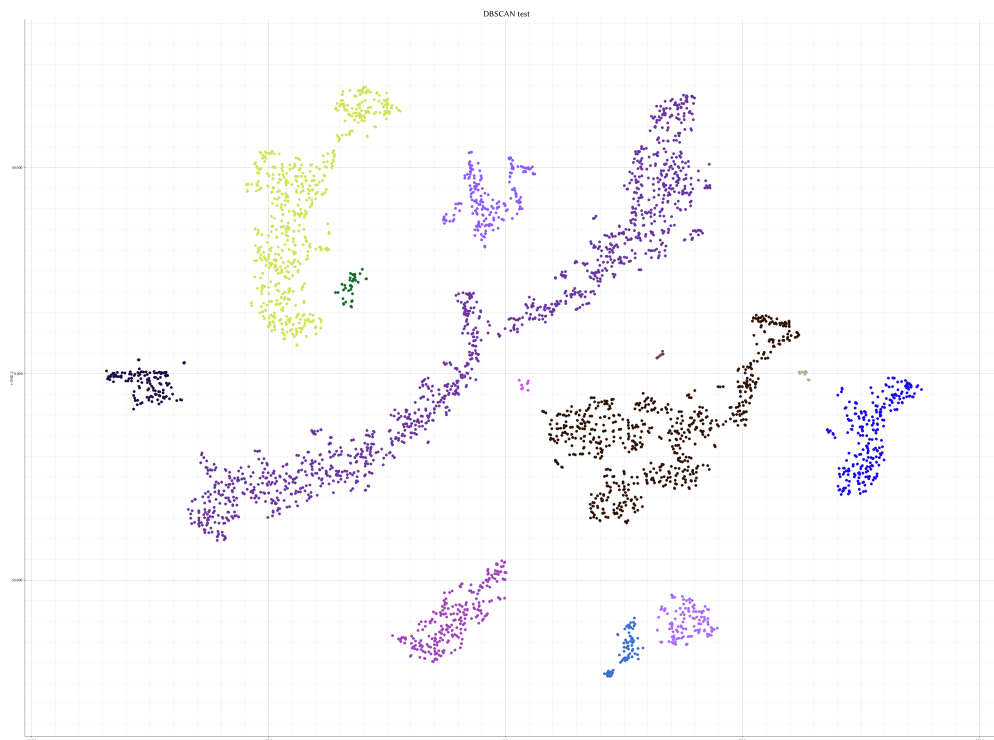


Рисунок 4.2 – Пример выполнения DBSCAN, параллельный алгоритм, 1 поток

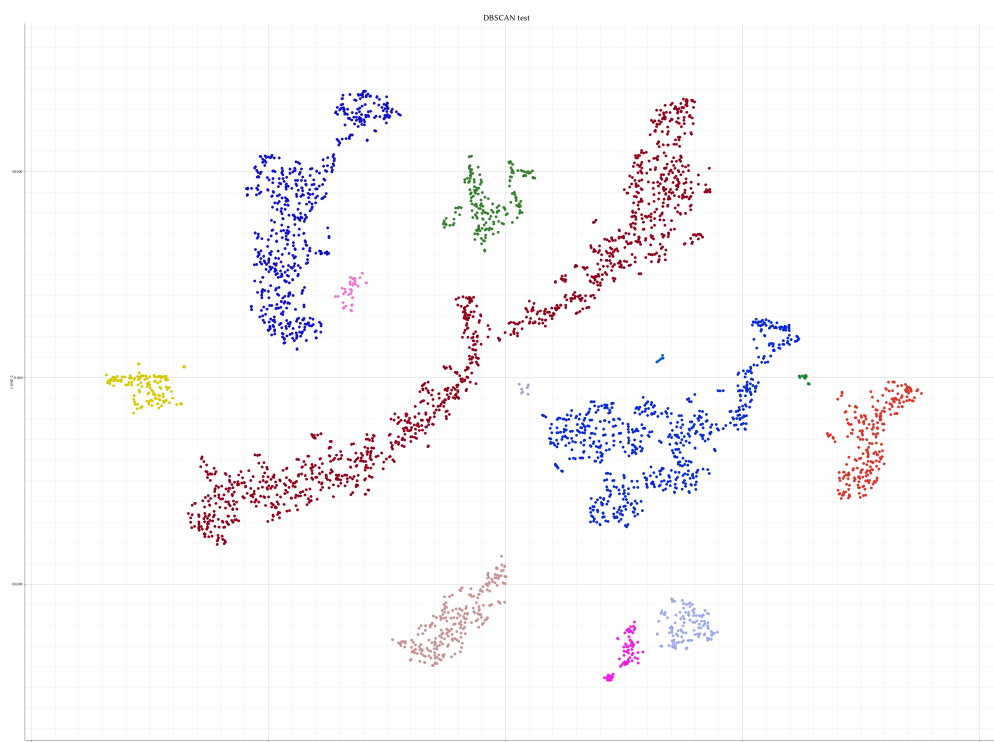


Рисунок 4.3 – Пример выполнения DBSCAN, параллельный алгоритм, 2 потока

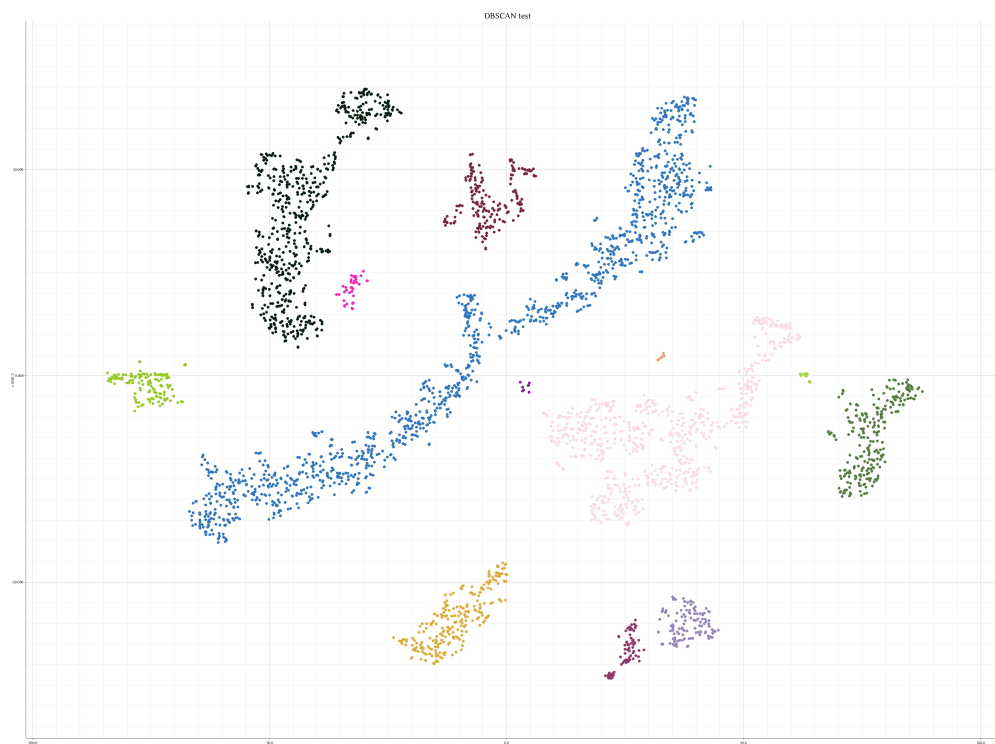


Рисунок 4.4 – Пример выполнения DBSCAN, параллельный алгоритм, 4 потока

4.2 Время выполнения алгоритмов

Алгоритмы тестировались при помощи инструментов замера времени предоставляемых библиотекой Criterion.rs[6]. Пример функции по замеру времени приведен в листинге 4.1. Количество повторов регулируется тестирующей системой самостоятельно, однако ввиду трудоемкости вычислений, количество повторов было ограничено до 25.

Листинг 4.1 – Пример функции замера времени

```
1 fn dbscan_bench(c: &mut Criterion)
2 {
3     let plot_config = PlotConfiguration::default().summary_scale(AxisScale::Linear);
4
5     let mut group = c.benchmark_group("DBSCAN");
6     group.plot_config(plot_config);
7     let path = select_file();
8     let cols = input_cols();
9
10    let dataset = Dataset::input_dataset(path.clone(), cols.clone(), b',');
11    group.sample_size(25);
12    let mut serial_model = Model::new(Arc::new(Mutex::new(dataset.clone())));
13    serial_model.set_eps(EPS);
14    serial_model.set_min_pts(MIN_PTS);
15    for size in 0..=6 {
16        // group.throughput(Throughput::Bytes(size_of_val(&*m1) as u64 + size_of_val(&*m2) as
17        // u64));
18        let size = 2usize.pow(size);
19        let mut parallel_model = ParallelModel::new(Arc::new(Mutex::new(dataset.clone())),
20            size);
21        parallel_model.set_eps(EPS);
22        parallel_model.set_min_pts(MIN_PTS);
23        group.bench_with_input(BenchmarkId::new("Parallel", size), &size, |b, _size| {
24            b.iter(|| {
25                parallel_model.run();
26                parallel_model.reset();
27                black_box();
28            });
29        });
30    }
31    group.finish();
32 }
```

График, показывающий время работы последовательного и параллельного алгоритмов в зависимости от количества потоков

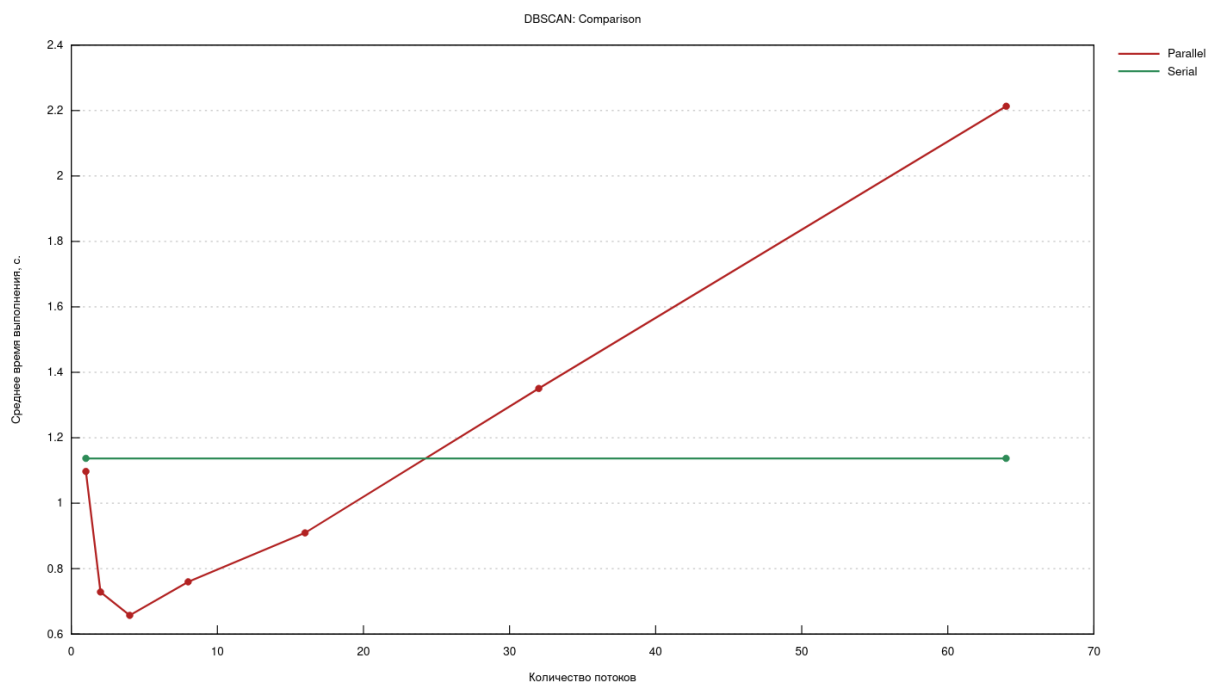


Рисунок 4.5 – Время работы DBSCAN

Как видно из полученных данных (на рис. 4.5), параллельная реализация работает быстрее, чем последовательная, при количестве рабочих потоков от 1 до 16. Более быстрая работа параллельного алгоритма при одном потоке по сравнению с последовательным алгоритмом может быть объяснена различными оптимизациями, которые были реализованы в ходе написания параллельного алгоритма. В частности, при параллельной реализации алгоритма мы не обязаны посещать все точки и просчитывать расстояние до каждой точки пространства — в случае, если ячейка содержит MinPt количество точек, то вся ячейка отмечается как часть кластера и нет необходимости в дополнительном просчете расстояния (гарантировано, что все точки внутри ячейки находятся на расстоянии меньше ϵ). Затраты, связанные с созданием поток(-а/-ов) при этом незначительны. Начиная с 8 рабочих потоков, процессорное время работы параллельной реализации растёт, и, при 32 рабочих потоках, время работы параллельного алгоритма превышает время работы последовательного. Такой рост связан с тем, что задействованный процессор не может эффективно обслуживать более 8 потоков, имея 8 логических ядер, что объясняет полученные результаты.

Вывод

В результате эксперимента было установлено, что параллельная реализация алгоритма DBSCAN работает быстрее, чем последовательная, при количестве рабочих потоков от 1 до 16 рабочих потоках. Это связано с тем, что использование более 8 рабочих потоков приводит к излишним временным затратам на диспетчеризацию потоков, поскольку на 1 логическое ядро приходится более 1 потока вычисления.

ЗАКЛЮЧЕНИЕ

В рамках данной лабораторной работы:

1. Были изучены и реализованы 2 версии алгоритма DBSCAN: последовательный и параллельный;
2. Был сделан сравнительный анализ алгоритмов на основе экспериментальных данных.
3. Был подготовлен отчёт по лабораторной работе, содержащий: особенности реализации многопоточной модели; временные характеристики предложенной реализации; результаты работы алгоритма;

Время работы последовательного алгоритма DBSCAN больше параллельного алгоритма, однако при неверном расчете количества потоков параллельный алгоритм значительно проигрывает своему аналогу.

Цели лабораторной работы по изучению и исследованию параллельных вычислений были достигнуты

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Автоматическая обработка текстов на естественном языке и компьютер-ная лингвистика : учеб. пособие / Большакова Е.И., Клышинский Э.С., Ландэ Д.В., Носков А.А., Пескова О.В., Ягунова Е.В. — М.: МИЭМ, 2011. — 272 с.
2. *Wang Y., Gu Y., Shun J.* Theoretically-Efficient and Practical Parallel DBSCAN. — 2019. — DOI: 10.48550/ARXIV.1912.06255. — Режим доступа: <https://arxiv.org/abs/1912.06255>.
3. Rust [Электронный ресурс]. — Дата обращения: 19.10.2022. Режим доступа: <https://www.rust-lang.org/>.
4. Arch Linux [Электронный ресурс]. — Дата обращения: 19.10.2022. Режим доступа: <https://archlinux.org/>.
5. Процессор Intel® Core™ i5-11320H [Электронный ресурс]. — Дата обращения: 19.10.2022. Режим доступа: <https://ark.intel.com/content/www/ru/ru/ark/products/217183/intel-core-i511320h-processor-8m-cache-up-to-4-50-ghz-with-ipu.html>.
6. Criterion [Электронный ресурс]. — Дата обращения: 19.10.2022. Режим доступа: <https://github.com/bheisler/criterion.rs>.