

A Fast DBSCAN Algorithm with Spark Implementation



Dianwei Han, Ankit Agrawal, Wei-keng Liao and Alok Choudhary

Abstract DBSCAN is a well-known clustering algorithm which is based on density and is able to identify arbitrary shaped clusters and eliminate noise data. Parallelization of DBSCAN is a challenging work because there is an inherent sequential data access order and based on MPI or OpenMP environments, there exist the issues of lack of fault-tolerance and there is no guarantee that workload is balanced. Moreover, programming with MPI requires data scientists to handle communication between nodes which is a big challenge. We present a new parallel DBSCAN algorithm using Spark. kd-tree technique is applied in our algorithm to reduce search time. More specifically, a novel merge approach is used so that no communication between executors is required while partial clusters are generated. Appropriate and efficient data structures are carefully used in our study: Using Queue to contain neighbors of the data point, and using Hashtable when checking the status of and processing the data points. Also other advanced data structures from Spark are applied to make our implementation more effective. We implement the algorithm in Java and evaluate its scalability by using different number of

D. Han (✉) · A. Agrawal (✉) · W. Liao (✉) · A. Choudhary (✉)
EECS Department, Northwestern University, Evanston, IL 60208, USA
e-mail: dianwei@eecs.northwestern.edu

A. Agrawal
e-mail: ankitag@eecs.northwestern.edu

W. Liao
e-mail: wkliao@eecs.northwestern.edu

A. Choudhary
e-mail: choudhar@eecs.northwestern.edu

processing cores. Our experiments demonstrate that the algorithm we propose scales up very well. Using data sets containing up to 1 million high-dimensional points, we show that our proposed algorithm achieves speedups up to 6 using 8 cores (10 k), 10 using 32 cores (100 k), and 137 using 512 cores (1 m). Another experiment using 10 k data points is conducted and the result shows that the algorithm with MapReduce achieves speedups to 1.3 using 2 cores, 2.0 using 4 cores, and 3.2 using 8 cores.

Keywords DBSCAN • Scalable data mining • Big data • Spark framework

1 Introduction

Clustering is a data mining approach that divides data into different categories that are meaningful, useful, or both [20]. Cluster analysis has been successfully applied to many fields: bioinformatics, machine learning, information retrieval, and statistics [20]. Well-known algorithms include K-means [13], BIRCH [24], WaveCluster [19], and DBSCAN [6]. Current clustering algorithms have been categorized into four types: partitioning based, hierarchy-based, grid-based, and density-based [6]. Density Based Spatial Clustering of Applications with Noise (DBSCAN) is a density based clustering algorithm [6].

Parallel DBSCAN has been implemented with MPI and OpenMP [4, 7, 15, 25]. Generally, an MPI implementation can obtain better performance but it requires the programmers to take care of implementation in detail, such as how to partition the data, how to deal with communication, synchronization, file location, and workload balancing. Besides parallelization with MPI, MapReduce-based approach is presented as well [7, 9, 14].

We propose a new distributed parallel algorithm with Spark that implements DBSCAN. A master-slave based approach is as follows. The algorithm first reads data from the Hadoop Distributed File System (HDFS) and forms Resilient Distributed Datasets (RDDs), transforming them into data points. Certainly, this process is done in Spark driver. It then sends the RDDs into multiple executors. Within each executor, partial clusters are generated and sent to driver at the end of *foreach* statement. Each executor just performs its computation without communicating with others. This way we avoid shuffle operations that are very expensive. So we place some additional points (SEEDs: the new term we introduce in our paper) in each partial cluster. After all the partial clusters are collected through shared variable accumulator, the algorithm identifies the clusters that are supposed to be merged by SEEDs. Merging is done in driver code too. In our new design and implementation, we use the power of shared variables of Spark framework:

broadcast and accumulator. Also, in order to shorten the search time for points' neighbors, we implement Java-based kd-tree [3] to reduce complexity from $O(n^2)$ to $O(n \log n)$. The experiments performed on a distributed-memory machine show that the proposed algorithm can obtain scalable performance.

The organization of the paper is as follows: In Sect. 2, we briefly give an overview of two frameworks based on big data: Map Reduce and Spark, and the basic idea of DBSCAN algorithm. Our proposed DBSCAN algorithm is introduced in Sect. 3. In Sect. 4, we present the parallel implementation with Spark. The experiments and the results are presented in Sect. 5, followed by some concluding remarks in Sect. 6.

2 Background

In this section, we first briefly review the basic idea of DBSCAN algorithm. And then we introduce two distributed computation frameworks that are very powerful and widely used in big data applications.

2.1 DBSCAN Algorithm

DBSCAN is a clustering algorithm proposed by Ester et al. [6]. And it has become one of the most common clustering algorithms because it is capable of discovering arbitrary shaped clusters and eliminating noise data [6]. The basic idea of this algorithm is finding all the *core points* and forming the clusters by clustering core points with all points (core or non-core) that are *reachable* from them. Essentially, DBSCAN algorithm is based on three basic definitions: core points, directly density-reachable, and density-reachable [25]. Given a data set D , of points.

eps-neighborhood of a point p is the neighborhood of $p \in D$ within a radius *eps*.

Definition 1 A point p is a *core point* if it has neighbors within a given radius (*eps*), and the number of neighbors is at least *minpts* (which is a threshold). In this case, the number of neighbors is called *density*.

Definition 2 A point y is *directly density-reachable* from x if y is within *eps-neighborhood* of x and x is a *core point*.

Definition 3 A point y is *density-reachable* from x if there is a chain of points p_1, p_2, \dots, p_n , with $p_1 = x, p_n = y$ and p_{i+1} is directly density-reachable from p_i for all $1 \leq i < n, p_i \in D$.

Algorithm 1 The DBSCAN algorithm

Input (eps , $minpts$, D)
Output (a set of clusters)

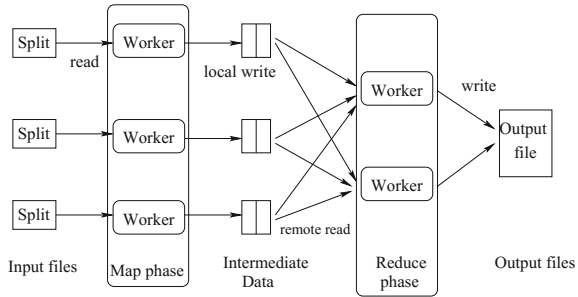
1. initialize all points as unvisited
2. **for** each unvisited point $p \in D$ **do**
3. mark p as visited
4. Let N be eps -neighborhood of p
5. **if** the size of $N < minpts$ points **then**
6. mark p as noise
7. **else**
8. create a new cluster C , and add p to C
9. **for** each point $p' \in N$
10. **if** p' is unvisited **then**
11. mark p' as visited
12. let N' be the eps -neighborhood of p'
13. **if** the size of N' is $\geq minpts$ **then**
14. add those points to N
15. **endif**
16. **endif**
17. **if** p' is not yet a member of any cluster
18. add p to C
19. **endif**
20. **endfor**
21. **endif**
22. **endfor**

The pseudocode of the DBSCAN algorithm is given in Algorithm 1 [8]. The algorithm starts with an arbitrary point $p \in D$ and checks its eps -neighborhood (Line 4). If the eps -neighborhood size is bigger than pre-defined number $minpts$, the code generates a new cluster C . The algorithm then retrieves all density reachable points from p in D , and add them to the cluster C (Line 8–20). Otherwise, if the eps -neighborhood contains less than $minpts$ points, then p is marked as noise (Line 6). The computational complexity of Algorithm is $O(n^2)$ where n is the number of data points. If we use spatial indexing, the complexity reduces to $O(n \log n)$ [3].

2.2 Two Powerful Frameworks Based on Big Data: MapReduce and Spark

In Hadoop version 1, MapReduce is the only data processing framework that is available for distributed computation. But in Hadoop version 2, based on Yarn (resource manager), MapReduce, Spark, and other data processing frameworks are

Fig. 1 An overview of data flow in MapReduce



available. MapReduce and Spark may share the same HDFS, but it should be pointed out that Spark jobs can be run with or without Yarn (Standalone mode).

- (1) *Map Reduce*: In big data domain, MapReduce is a simple but powerful framework which makes programmer easily implement parallel processing. It is based on Hadoop Distributed File System (HDFS), which allows programmers to focus mainly on the problem itself instead of the low level implementation details. Figure 1 tells us about how this programming model works. MAP workers read data from HDFS and process the data based on the business logic and then write intermediate data to local disk for sorting and shuffling process. It is also in the form of key-value pair. After a reduce worker is notified by master, it uses remote procedure call to read data from local disk of MAP workers, and then sorts data so that all occurrences of the same key are grouped together. The output of reduce function will be appended to final output files (generally HDFS).

Compared with the other distributed computation frame-work, MapReduce has the following advantages:

- **Extremely Scalable.** It does not require the support from centralized RAID-based SAN or NAS storage systems. Every node has its own local hard-drives. The nodes are loosely coupled and connected with standard network devices. So adding and removing nodes to a cluster becomes very easy and convenient, and has no impact to running MapReduce jobs [17].
- **Highly Parallel and Abstracted.** Based on the frame-work's principle, programmers do not have to take care of low level implementation details such as message transferring between master and workers, file location, and workload balancing. They only need focus on the problem itself. One of the major contributions of MapReduce is that it supports parallelization automatically. The programmers only need to implement `map()` method of Mapper class and `reduce()` method of Reducer class and the framework will do the rest. However, for

complicated job, the programmers still need to figure out the number of Mappers and how to split the input data.

- Highly Reliable and Fault-tolerant. From the data source perspective, HDFS uses the replication strategy to handle data source reliability. A single process failure in MPI will cause the whole job to fail. In MapReduce framework, another task will be automatically launched if one task fails and the job will continue running. This feature is especially useful and important for long-running jobs.
- (2) *Spark*: At a high level, a running Spark application has one driver process talking to many executor processes, sending them work to do and collecting the results of that work. The first thing a Spark program must do is to create a `SparkContext` object in driver code, which tells Spark how to access a cluster. Then it reads one file or multiple files in HDFS and processes them as Distributed Datasets (RDD), which is a collection of elements partitioned across the nodes and can be operated on in parallel. RDD is the main abstraction Spark provides, and RDDs can be created from a file in the Hadoop file system or by transforming other RDDs. We want to point out that Spark can use not only its own APIs to read data but also Hadoop API to read data (`newAPIHadoopFile` method in this case). `TaskScheduler` launches tasks to executors via Resource manager, which in this case, is YARN. After executors complete their tasks, they will send the results back to the driver (see Fig. 2) (if it is the final RDD of an action such as `count()`) [12], or write output to external storage. Spark framework captures all the important features that MapReduce have. In addition, it has the following new features.
- In-memory computations. In Spark, Resilient Distributed Datasets (RDDs) are the first abstraction that allows programmers to perform in-memory computations on large clusters. RDDs are motivated by two types of applications that MapReduce handle inefficiently: iterative algorithms and interactive data mining [22]. Figure 1 depicts that MapReduce frameworks does not fit iterative algorithms. In order to use MapReduce model to tackle iterative algorithms, many rounds of map-reduce executions will be performed which is not very efficient because map's intermediate results should be written to local disks and then they are remotely read to reduce workers, and disk I/O operations are very expensive in this case. In Spark, the benefit of keeping everything in memory is the ability to perform iterative computations at blazing fast speeds.

- Supporting Streaming data, complex analytics, and real time analysis. MapReduce offers a very simple but powerful programming model that are efficient for data-intensive algorithms [11]. But we can not use MapReduce to perform real time analysis and implementing complex graph based algorithms in an efficient manner.
- Fast fault recovery. In MapReduce old version, if the JobTracker does not receive any heartbeat from a TaskTracker for a specified period of time, the JobTracker understands that the worker associated to that TaskTracker has failed. When this situation happens, the JobTracker needs to reschedule all pending and in-progress tasks to another TaskTracker, because the intermediate data belonging to the failed TaskTracker may not be available anymore [21]. After hadoop-0.21, checkpointing was added where JobTracker records its progress in a file. When a JobTracker starts, it can restart work from where it left off. MapReduce uses replication strategy to handle fault recovery. On the other hand, Spark reconstructs RDDs via lineage to handle this issue. Compared to the replication method, which consumes more memory, reconstruction of RDDs takes shorter time [23].

Even though spark is very efficient, offers parallelization automatically, we still need to put much effort to avoid shuffle operation. So in our implementation of DBSCAN we avoid all-to-all communication.

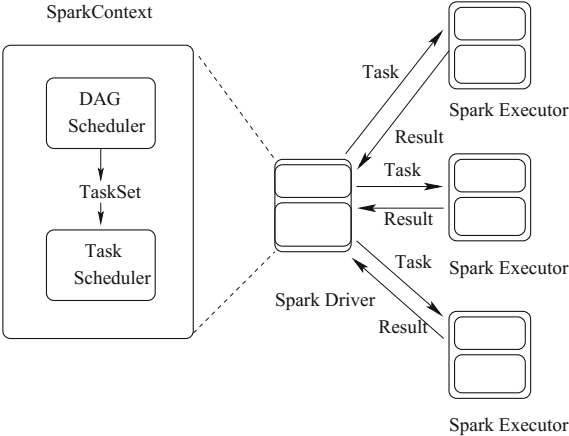


Fig. 2 An overview of data flow in spark

3 Novel DBSCAN with Spark Implementation

To our best knowledge, there are many DBSCAN implementations with Hadoop's MapReduce [9, 14, 17]. But very few people implement DBSCAN with Spark because the programmers need to design a new algorithm to avoid shuffle operations to make their parallelization more efficient. For example, after one data point's state is updated in one executor we need to spread this updating across the cluster. So this will introduce shuffle operations which are very expensive in Spark. Let us take a look at the pseudocode of our new DBSCAN's algorithm.

3.1 DBSCAN Algorithm with Spark

The pseudocode of the DBSCAN algorithm with Spark implementation is given in Algorithm 2. The algorithm starts with the code in Spark driver, which reads data, generates RDDs and transforms them into appropriate RDDs (Line 1, Line 2, and Line 3). The code in Spark executor is in Lines 4 through 32. After comparing with Algorithm 1, we can see that two places are new: Line 15 and Lines 29 through 31. We assume each executor only deals with the points that belong to it. Otherwise, there would be a lot of overlap of computation between different executors. Placing SEEDs is in Line 15. The detailed description regarding it will be given in next Section. The partial clusters are sent back to driver right before the executor finishes its task by accumulator, which also will be explained in detail in the next Section. This implementation is meant for ensuring that merging process will not be started until all the executors finish their tasks. Lines 33 through 34 perform merging partial clusters and produce the final global clusters (see Algorithm 4). So the code [1–3] is run in Driver mode, code [4–32] is run in Executor mode, and code [33–34] is run in Driver mode.

Algorithm 2 DBSCAN algorithm with Spark

Input(*eps*, *minpts*, *D*)**Output** (a set of clusters)

```

1.  read an input file from HDFS and generate RDDs from the read data
2.  transform the existing RDDs into appropriate RDDs with Point type
3.  distribute those RDDs into executors
4.  foreach (closure's start)
5.      if point p is not in hashtable then
6.          get the neighbors of p using eps and kdtree
7.          push all appropriate neighbors into Queue N
8.          if the size of N < minpts points then
9.              mark p as noise
10.         else
11.             create a new cluster C, and add p to C
12.             while N is not empty
13.                 let p' be the removed point from N
14.                 put the index of p' into the hashtable
15.                 place SEEDs processing
16.                 if p' is unvisited then
17.                     mark p' as visited
18.                     let N' be the eps-neighborhood of p'
19.                     if the size of N' is  $\geq$  minpts then
20.                         add appropriate points to N
21.                     endif
22.                 endif
23.                 if p' is not yet a member of any cluster
24.                     add p' to C
25.                 endif
26.             endwhile
27.         endif
28.     endif
29.     if current point is the last one in closure then
30.         send partial clusters to driver through accumulator
31.     endif
32. endforeach
33. analyze partial clusters based on the placed SEEDs
34. search for all partial clusters and merge them if necessary

```

3.2 Two Important Data Structures Affecting Performance

Using Java as the programming language in our implementation, we need to consider using the appropriate data structures for efficiency. Here, two data structures Hashtable and Queue are discussed.

If we take a look at Line 14, this operation should be *put(key, value)*, which is usually $O(1 + n/K)$ where K is the hash table size. If K is large enough, the result is effectively $O(1)$. Method *containsKey(key)* is performed in Line 5, Line 7, and Line 20. Again, under normal circumstances, it is $O(1)$. The add operations on Queue are performed in Line 7 and Line 20, and remove operation on Queue is performed in Line 13. The number of add operations should be the same as the number of remove operations according to the condition in Line 12 (while loop will not terminate until it is empty). Among LinkedList, ArrayList, and Vector, the best performance on both add and remove operations is obtained using LinkedList. In our code, we thus use LinkedList to implement Queue.

4 Novel Techniques in Parallel DBSCAN with Spark

In this Section, we will present the implementation details of our parallel DBSCAN algorithm with Spark. The pseudocode of algorithms is given in the first part. Then we analyze the time complexity of the whole algorithm.

Algorithm 3 Placing SEEDs in Executors

```

1.  identify the current partition of this executor as par_A
2.  initialize the place_flg for all the partitions
3.  while  $N$  is not empty
4.      let  $p'$  be the removed point from  $N$ 
5.      put the index of  $p'$  into the hashtable
6.      for  $j = 0 \dots$  all the possible partitions
7.          if  $p'$  is in par_A
8.              let continue_flg be true
9.              break
10.         else
11.             if place one seed already
12.                 let continue_flg be false
13.                 break
14.             else
15.                 place_flg = 1
16.                 let continue_flg be true
17.                 break
18.             endif
19.         endif
20.     endfor
21.     if continue_flg = false
22.         continue
23.     endif
24.     if place_flg is 1
25.         place a seed
26.     endif
27. endwhile

```

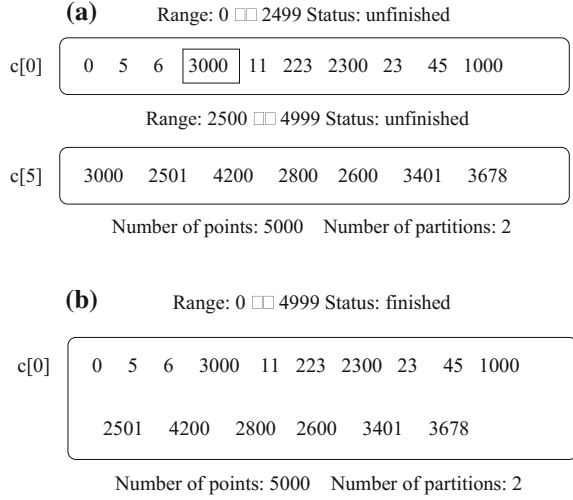
4.1 New Clustering Algorithm Without Communication Between Executors

We need to update data points' state by map function if we apply the traditional method, and then propagate this update to other executors. However, that implementation will introduce a shuffle operation in order to make this update visible by other executors. Here, we propose a novel clustering algorithm to get around the shuffle operation. After data points have been partitioned to each executor, we just let each executor compute the partial clusters locally for data points that are assigned to this executor. The merging process is deferred until all the partial clusters have been sent back to the driver. This new design, however, introduces new challenges: how to create the partial clusters in executors so that they can be merged in the driver? And how to identify those partial clusters which are supposed to be merged into one cluster? The pseudocode of algorithms and an example are given as follows.

Algorithm 3 gives the basic idea of our design. In order to avoid overlap of computation of partial clusters, we would let individual executors only deal with the points that belong to this partition so that the executors do not have to communicate to spread points' updated states across the clusters. However, we could not merge the partial clusters into the global clusters after all the partial clusters are collected in driver because there are no global states of these partial clusters. Therefore, we come up with a new idea, using SEEDs, which are points that do not belong to the current partition. And these SEEDs serve as something like markers so that we can easily identify the outer master partial clusters by using them and merge them into a bigger cluster. The SEEDs are not related to the locations. If the current point's index is beyond the range of current partition it is taken as a SEED. So the main goal on executor side is to place SEEDs, and on driver side, we find out SEEDs and identify master partial clusters and merge them.

Before moving on to the algorithm of digging out SEEDs from partial clusters in Spark driver, we would like to use an example to display how to identify SEEDs and search for master partial clusters. Figure 3a shows that there are 2 partial clusters from 2 partitions. SEEDs are those points whose indexes are beyond the partition's range. For example, for $C[0]$, its range is from 0 to 2499. So the point whose index is greater than 2499 is 3000. Then the algorithm will identify the master partial clusters. Obviously, for 3000, the master partial cluster is $C[5]$ because it contains 3000 and 3000 is a regular element in this cluster. When we merge two partial clusters we need to remove duplicate elements. Figure 3b show the resulting cluster $C[0]$.

Fig. 3 An example showing the proposed merging cluster algorithm at different stages. **a** There are two partitions and two partial clusters. Integers in squares are SEEDs. **b** After C[0] merges C[5], C[0] status is updated as “finished” from “unfinished”



In Spark driver, Algorithm 4 shows how to use SEEDs to merge partial clusters into global clusters. First of all, it identifies the SEEDs by comparing elements with its range. In general, the number of SEEDs should be equal to or greater than the number of partitions. So we obtain an array of seeds (see Line 3). Lines from 4 through 8 form a for loop, which finds the master cluster that contains the seed as a regular element, then merges the two clusters, and finally, updates the status of master cluster. When the for loop terminates the status of current cluster is updated from ‘unfinished’ to ‘finished’.

Algorithm 4 Using SEEDs and merge partial clusters in Driver

```

1.  for i = 0.. all partial clusters
2.    if the status of current partial cluster is ‘unfinished’
3.      seed = identify seeds from current partial cluster
4.      for j = 0..seed.size()
5.        rtn_index = find master partial cluster index
6.        merge current with master cluster
7.        update the status of master cluster to ‘finished’
8.      endfor
9.    update the status of current cluster to ‘finished’
10.  endif
11. endfor

```

4.2 Time Complexity Analysis

We define some related notations as follows:

n the number of data points;
 p the number of partitions;

m	the number of partial clusters;
K	the maximum size of partial clusters;
$t_{straggling}$	the average wait time for framework to allow all stragglers to finish.
T_s	the average time complexity of the sequential algorithm;
T_p	the average time complexity of the parallel algorithm;
$Save$	the average speed-up.

Basically, there are three parts in our algorithm.

In the first part, the driver reads data points from HDFS and transforms the data points into appropriate form that can be processed in executors and constructs the *kd-tree*. The time for this phase includes reading file, transforming RDDs, and building *kd-tree*. We assume we use Δ for the first two items. For *kd-tree* construction, we use $O(n \log n)$ [10]. So summing them up, we use $\Delta + O(n * \log n)$.

In the second part, the local partial clusters are generated in executors. Basically, in the best case, searching a point from a balanced *kd-tree* takes $O(\log n)$ time. In the worst case, the time could be n . Some researchers have reported that (near neighbor) range search's upper bound is $O(nl - l/d + k)$ [10]. So we use V to represent the search time, which is between $\log n$ and $n - l/d + k$; If we use parallel processing, we need to add the time for SEEDs placement part. Let us assume an additional $O(m * V)$ time is added. So in parallel processing, we would spend $O((n/p * V) + (m * V)) + t_{straggling}$ time in our case.

In the last part, after all the partial clusters have been sent back from executors to the driver, the driver merges them and produces the global clusters. Based on our Algorithm 4, the search operations takes $O(n)$ time at most if we check each element in the partial clusters. For merging phase, it takes Km times which is less than n . So we use $O(n + Km)$ time.

To sum up: $T_s = O(\Delta + n * \log n + n * V + n + Km)$.

$$T_p = O(\Delta + n * \log n + (n/p) * V + m * V + t_{straggling} + n + Km)$$

$$Save = T_s / T_p.$$

5 Experiments and Analysis

A series of experimental tests are conducted to evaluate the effectiveness and efficiency of our DBSCAN algorithm with Spark and MapReduce's implementations. We need to note that all parallel executions generate the same result as the serial execution. The dimension of data is relevant to the computational cost of querying the *kd-tree*. We do not perform tests based on varying number of attributes because we focus on Spark implementation instead of *kd-tree* implementation in our work. The tests are done on different sizes of data points with multiple dimensions. Our experimental results have been reported in terms of the CPU times.

After comparing with the results from Patwary et al. [15], we find that our results match them so we do not list the accuracy in our paper.

5.1 Experimental Setup

To perform the experiment for our DBSCAN’s parallel implementation with Spark, we use Edison (operated by Lawrence Berkeley National Laboratory and the Department of Energy Office of Science), a Cray XC30 distributed memory parallel computer. It has 5576 compute nodes, 133,824 cores in total. Each node has two 12-core Intel “Ivy Bridge” processors at 2.4 GHz and 64 GB DDR3 1866 MHz memory. Each core has its own L1 and L2 caches, with 64 KB (32 KB instruction cache, 32 KB data) and 256 KB, respectively; A 30 MB L3 cache shared between 12 cores on the “Ivy Bridge” processor [5]. The algorithms have been implemented in Java (1.7) using the Spark (1.5) and Hadoop (2.4).

Our testbed consists of 5 datasets, which are divided two groups: (c10 k, c100 k), and (r10 k, r100 k, r1 m). Both groups of datasets (*synthetic-cluster*) have been generated synthetically using the IBM synthetic data generator [1, 16]. Table 1 lists the properties of our test data.

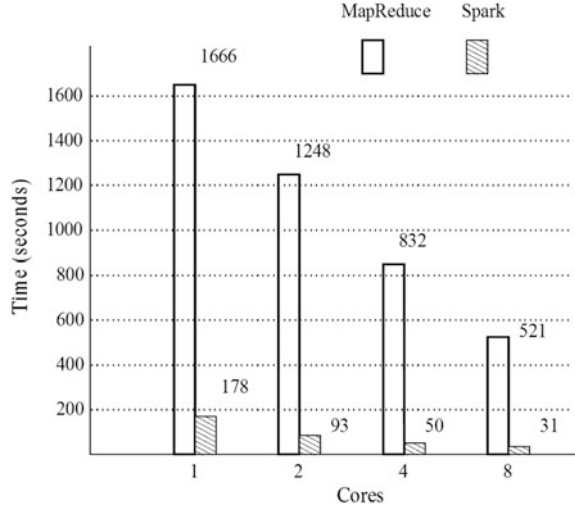
5.2 Comparison of the Time Taken by MapReduce and Spark

As we are not able to get source code from the other research teams [7, 9, 14], we have implemented our own DBSCAN with MapReduce approach. From Fig. 4, it is seen that 9–16 times faster performance is obtained from Spark than MapReduce. Due to the length of time taken by MapReduce, we have not conducted further tests on medium scale and large scale data sets.

Table 1 Properties of test data

Name	Points	d	eps	Minpts
c10 k	10,000	10	25	5
c100 k	102,400	10	25	5
r10 k	10,000	10	25	5
r100 k	102,400	10	25	5
r1 m	1,024,000	10	25	5

Fig. 4 Time used by MapReduce and spark.
Number of points: 10,000,
dimension: 10, eps, 25.0,
minPts: 5



5.3 Comparison of the Time Spent in Driver and in Executors

In this part, we discuss the time taken in our program. Figure 5a–d shows the time taken between executors and driver according to our experiments. Based on the Algorithm 2, we expect to see more time will be spent in driver with the number of partial clusters increasing. Let us take a look at Fig. 5a first. When we use more cores (1–8) to run our program, we see the number of partial clusters becomes bigger (10–392), but the time spent in driver does not change very much. That is because the data set is too small. Take a look at Fig. 5c, d, their patterns are exactly the same. When using more cores (4–32), more partial clusters are produced (from 720 to 9279), and the time spent in driver gradually becomes more. This is consistent with our analysis on the time complexity that we conduct in Sect. 4, where when the number of partial clusters m increases, the time $n + Km$ becomes large as well. Figure 5b follows the complexity analysis as well.

5.4 Scalability of Parallel DBSCAN with Spark

Before we discuss the scalability of our algorithm, we need to mention that for large data sets (≥ 1 million data points), we use *kd-tree* with pruning branches to shorten search time.

The speedup obtained by our DBSCAN algorithm with Spark is given in Fig. 6. The left column in Fig. 6 shows the speedup considering only the computation in executors while the right column shows the results considering the computation in

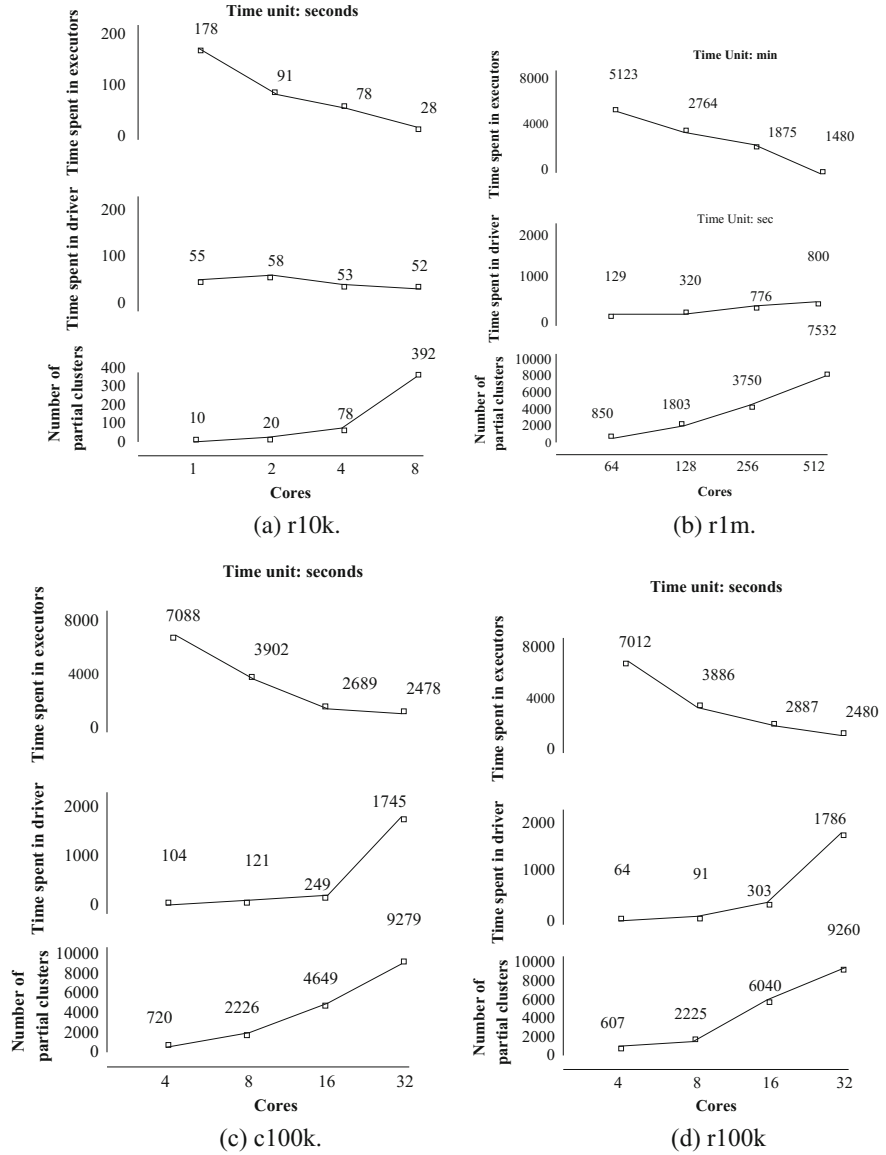


Fig. 5 The time distribution between driver and executors

executors and driver. It is obvious that the local computation in executors scales better than the whole computation since their computations are independent. For 10 k data sets, we obtain speedup up to 1.9, 3.6, and 6.2 respectively using 2, 4, and 8 cores. For 100 k data sets, speedup up to 3.3, 6.0, 8.8, and 10.2 respectively using

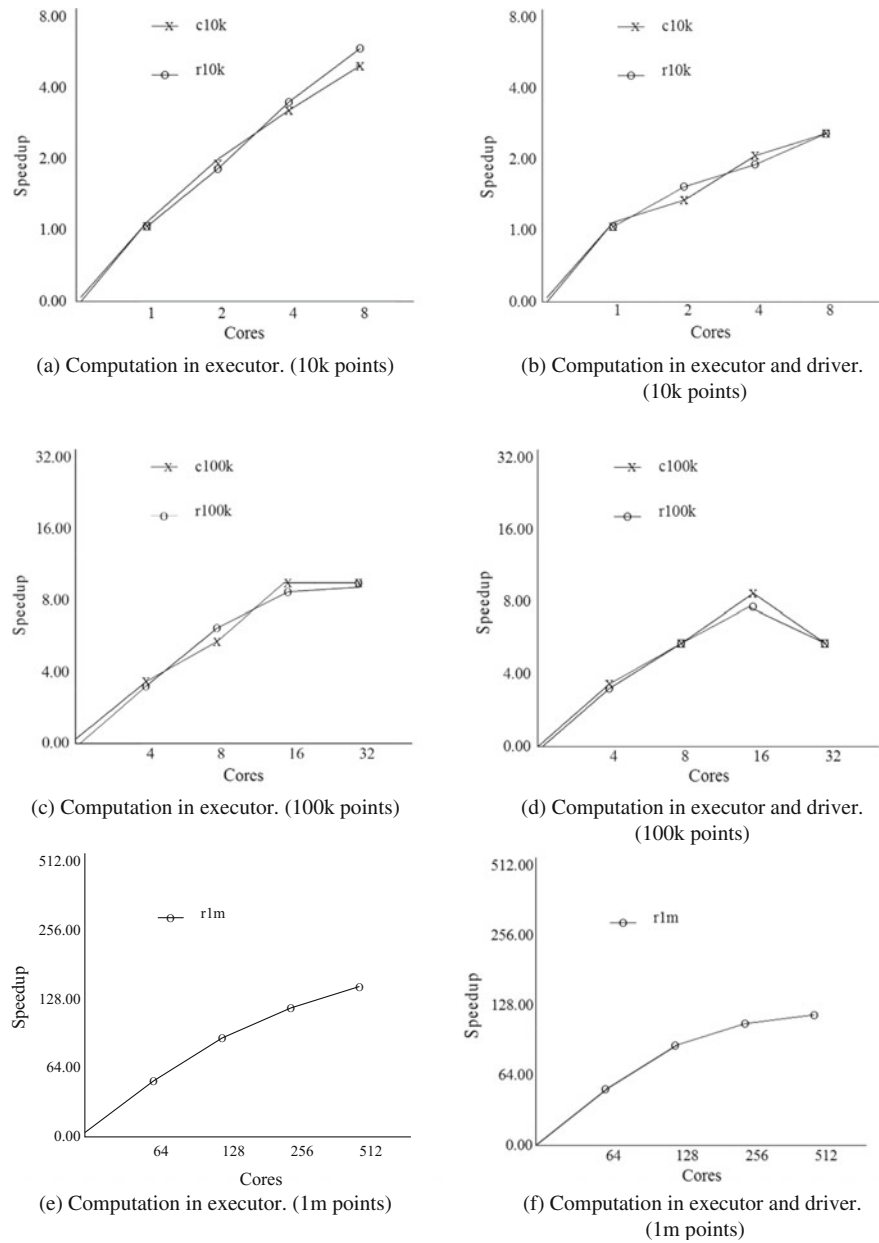


Fig. 6 Speedup of DBSCAN algorithm with spark. Left side: time spent in executor. Right side: time spent in driver and executor

4, 8, 16, and 32 cores. For 1 m data set, speedup up to 58, 83, 110, and 137 respectively using 64, 128, 256, and 512 cores.

Take a look at right column, Fig. 6b, d, f show the speedup when total time is considered. The curves seem more flat compared with the ones in left column. For 10 k data sets, because the total time is less, the merging time is not significant. For 100 k data sets, more partial clusters are collected in driver. When using 4, 8, and 16 cores, the local computation time still dominates the total time, so speedup does not change very much. When using 32 cores, 9279 partial clusters are generated in executors and collected in driver. So the speedup drops to 5.6 from 10.2.

For r1 m, we use pruning branches technique, and thus the neighbor size of each point is decreased. Also we filter out those partial clusters whose size is too small, and their removal does not impact the accuracy significantly. Therefore, the speedup of total time does not change a lot compared with local computation.

6 Conclusions

DBSCAN algorithm has been very powerful and popular because it is able to identify arbitrary shaped clusters as well as handle noisy data. However, parallelization of DBSCAN based on MPI and OpenMP suffers from lack of fault-tolerance. Moreover, in order to implement parallelization with MPI or OpenMP, data scientists need to take care of implementation in detail, such as handling communication, dealing with synchronization, and so forth, which can pose a challenge for many users. In this paper, we proposed a new Parallel DBSCAN algorithm with Spark, which avoids the communication between executors and thus leads to a better scalable performance. The results of these experiments demonstrate that our new DBSCAN algorithm with Spark is scalable and outperforms the implementation based on MapReduce by a factor of more than 10 in terms of efficiency. In the future, we would try to apply partitioning strategy with Spark implementation and try to use larger datasets in our study.

Acknowledgements This work is supported in part by the following grants: NSF awards CCF-1409601, IIS-1343639, and CCF-1029166; DOE awards DESC0007456 and DE-SC0014330; AFOSR award FA9550-12-1-0458; NIST award 70NANB14H012. This research used Edison Cray XC30 computer of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.

References

1. Agrawal, R., & Srikant, R. (1994). Quest synthetic data generator, *IBM Almaden Research Center*.
2. Beckmann, N., et al. (1990). The r^* -tree: An efficient and robust access method for points and rectangles. In: *Proceedings of the 1990 ACM SIGMOD International Conference on Management of Data* (Vol. 19, no. 2, pp. 323–331).
3. Bentley, J. (1975). Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18(9), 509–517.
4. Brecheisen, S., et al. (2006). Parallel density-based clustering of complex objects. *Advances in Knowledge Discovery and Data Mining*, pp. 179–188.
5. DOE Office of Science (2015, September 17). Edison Configuration (Online). <https://www.neresc.gov/users/computational-systems/edison/configuration/>.
6. Ester, M., et al. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In: *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining* (Vol. 1996, pp. 226–231). AAAI Press.
7. Fu, Y., et al. (2011). Research on parallel DBSCAN algorithm design based on mapreduce. *Advanced Materials Research* 301, 1133–1138.
8. Han, J., et al. (2011). *Data mining: Concepts and Techniques*. Morgan Kaufmann.
9. He, Y., et al. (2014). MR-DBSCAN: A scalable mapreduce-based DBSCAN algorithm for heavily skewed data. *Frontiers of Computer Science*, 8(1), 83–99.
10. Kakde, H. M. (2005, August 25). *Range Searching using Kd Tree* (Online). <http://www.cs.utah.edu/lifeifei/cs6931/kdtree.pdf>.
11. Kang, S. J., et al. (2015). Performance comparison of OpenMP, MPI, and MapReduce in practical problems. *Advances in Multimedia* 2015.
12. Karau, H., et al. (2015). *Learning Spark: Lightning-fast Data Analysis*. O'Reilly Media.
13. MacQueen, J., et al. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability* (Vol. 1, pp. 281–297). USA.
14. Noticewala, M., & Vaghela, D. (2014). MR-IDBSCAN: Efficient parallel incremental DBSCAN algorithm using mapreduce. *International Journal of Computer Applications* 93(4), 13–17.
15. Patwary, M. M. A., et al. (2012). A new scalable parallel DBSCAN algorithm using the disjoint-set data structure. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC 2012*, pp. 62:1–62:11. IEEE Computer Society Press.
16. Pisharath, J., et al. (2010). *NU-MineBench 3.0*. Technical Report CUCIS-2005-08-01, Northwestern University (Technical Report).
17. Sakr, S., & Gaber, M. M. (2014). *Large Scale and Big Data: Processing and Management*. CRC Press.
18. Spark, A. (2015). *Spark Programming Guide* (Online). <http://spark.apache.org/docs/latest/programming-guide.html>.
19. Sheikholeslami, G., et al. (2000). WaveCluster: A wavelet based clustering approach for spatial data in very large databases. *The VLDB Journal*, 8(3), 289–304.
20. Tan, P., et al. (2005). *Introduction to Data Mining*. Pearson.
21. White, T. (2011). *Hadoop: The Definitive Guide*. O'Reilly Media.
22. Zaharia, M., et al. (2012) Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation* (pp. 2–2). USENIX Association.

23. Zaharia, M. (2014). *An Architecture for Fast and General Data Processing on Large Clusters*. Technical Report UCB/EECS-2014-12, University of California, Berkeley (Technical Report).
24. Zhang, T., et al. (1996). BIRCH: An efficient data clustering method for very large databases. In *ACM SIGMOD Record* (Vol. 25, Issue. 2, pp. 103–114). ACM.
25. Zhou, et al. (2000). Approaches for scaling DBSCAN algorithm to large spatial databases. *Journal of Computer Science and Technology*, 15(6), 509–526.