

Aalborg University

Department of Computer Science

Selma Lagerlöfs Vej 300

9220 Aalborg Ø

Telephone 99 40 99 40

<http://www.cs.aau.dk>

Title:

Efficient Skyline Computation
for Large Volume Data in
MapReduce Utilising Multiple
Reducers

Theme:

Database Technology

Semester:

Software Engineering

10th semester

Feb. 4th, 2013 - June 7th, 2013

Group:

sw106f13

Members:

Kasper Mullesgaard

Jens Laurits Pedersen

Supervisor:

Hua Lu

Copies: Four

Paper - pages: 33

Appendices: 1

CD-ROM: 1

Total pages: 35

Abstract:

A skyline query is useful for extracting a complete set of interesting tuples from a large data set according to multiple criteria. The sizes of data sets are constantly increasing and the architecture of backends have switched from single node environments to cluster oriented setups. Previous work has presented ways to run the skyline query in these setups using the MapReduce framework, but the parallel possibilities are not taken advantage of since a significant part of the query is always run serially. In this paper, we propose the novel algorithm Grid Partitioning Multiple Reducers Skyline (GPMRS) that runs the entire query in parallel. This means that GPMRS scales well for large data sets and large clusters. We demonstrate this using experiments showing that GPMRS runs several times faster than the alternatives for large data sets with high skyline percentages.

Summary

The purpose of this master thesis project is to develop and implement an efficient skyline query algorithm in the MapReduce framework. The project is documented in the form of a scientific article and the focus is on utilising multiple reducers to process large data sets.

Two algorithms are proposed, one being an extension of the other. The first one is called MR-GPSRS and it was primarily developed during last semester. It is an efficient skyline query processing MapReduce algorithm. The second algorithm is called MR-GPMRS. It is an extension to MR-GPSRS, and unlike MR-GPSRS, it utilises multiple reducers to compute the global skyline.

To evaluate MR-GPSRS and MR-GPMRS, both algorithms are compared to two other algorithms from the literature, MR-Angle and MR-BNL. The experiments show that MR-GPSRS performs better than MR-Angle and MR-BNL in all tests, and that MR-GPMRS is better than MR-GPSRS when the skyline percentage of the data set is high.

The proposed algorithms perform better than the algorithms from the literature because of two reasons. The first being that by utilising grid partitioning they are able to prune a large part of the domination checks between tuples. The second being that MR-GPMRS scales better with the skyline percentage as it allows the entire skyline computation to be performed parallelly. In addition, a cost estimation is made that provides a worst case number of partition comparisons in MR-GPMRS.

To conclude the paper, some future research directions are suggested. The evaluation conducted in the paper was performed on a relatively small cluster, and in order to properly test the algorithms performance on big data, a bigger cluster is needed. Another direction is to find methods for further reducing the communications cost between the mappers and reducers while load balancing the reducers.

Kasper Mullesgaard

Jens Laurits Pedersen

Date: June 7th, 2013

Contents

1	Project Overview	7
1.1	1 st Semester	7
1.2	2 nd Semester	8
1.3	Reflection	9
2	Scientific Article	11
3	Implementation	25
3.1	MR-GPSRS and MR-GPMRS	25
3.1.1	Bitstring Generation	25
3.1.2	Skyline Computation	26
3.2	MR-BNL	29
3.2.1	Phase 1	29
3.2.2	Phase 2	31
3.3	MR-ANGLE	32
3.3.1	Phase 1	32
3.3.2	Phase 2	33
A	CD-ROM	35

Chapter 1

Project Overview

This project is the culmination of the work across two semesters. The 1st semester lasted from September 2012 to January 2013, and the 2nd from February 2013 to June 2013. This paper documents the work in this semester. This paper consists of three chapters. This chapter contains an overview of the entire project: A summary of the two semesters and a project reflection. Chapter 2 contains the scientific article that is the main contribution of this paper, and Chapter 3 contains a description of how the algorithms associated with this project have been implemented.

1.1 1st Semester

The initial problem statement was to develop an efficient skyline query processing algorithm for MapReduce. The current research into this area has room for improvement as all previous work relies on having the final skyline computed at a single node. The final goal was therefore to develop an algorithm where the final skyline is computed parallelly across multiple nodes.

The first goal in the previous semester was to establish a foundation for the skyline query processing in MapReduce. This included establishing familiarity with the Hadoop MapReduce framework and how it implements the MapReduce framework.

The second goal was to implement a skyline query processing algorithm. The initial idea was to solve utilise multiple reducers by copying the data for calculating the final skyline to multiple nodes, each node having responsibility for a part of the data set. In this way, multiple nodes were able to calculate the final skyline parallelly. This method, however, was very naive and even with the multiple reducers, it lacked the efficiency to compete with the current research. The reason was a lack of the ability to prune any domination checks between tuples.

The project focus then moved on to improve pruning power. In skyline

query processing, grid partitioning is an efficient way of pruning domination checks between points. Implementing grid partitioning requires that the mappers and reducers are aware of which partitions are empty and non-empty globally. This is problematic since mappers and reducers cannot communicate with each other. The solution to this was to generate a bitstring that represented the empty and non-empty status of partitions across the entire data set, in a MapReduce job before the skyline computation. The MapReduce job that computed the skyline was then initialized with the generated bitstring, allowing the mappers to make pruning decisions where the entire data set is taken into account. The algorithm that resulted from this is called MR-GPSRS.

The next goal was to find a new way to utilise multiple reducers that made effective use of the grid partitioning scheme. MR-GPSRS only uses a single reducer, since the original method for using multiple reducers found in the start of the semester was too naive to continue using, as simply copying the data set across all reducers caused an unnecessary high communication cost. Computing the global skyline efficiently using multiple reducers was not solved by the end of the 1st semester. A method was found for the mappers to send individually processable data to multiple reducers, but the method caused the communication cost to become extremely high. This was named the multiple reducers problem. Another problem with the algorithm at the end of the 1st semester was that the Partitions Per Dimension (PPD) had to be manually set. The most efficient PPD could only be determined by having intricate knowledge of the data set, or by running the algorithm multiple times with different values. This was named the tuning problem.

1.2 2nd Semester

The previously identified problems from the 1st semester was the focus of the work during the 2nd semester. The multiple reducer problem was the main focus since solving it would be a novel and valuable contribution to skyline query processing in MapReduce. A solution was found, and further time went into refining and improving certain aspects of it. The solution to the problem was a method that allowed for identifying groups of partitions that could be processed independently at the reducers. These groups are then merged into a number of groups equal to the number of reducers, and the groups are sent to the reducers where they each compute a part of the global skyline. This method was implemented as a new algorithm MR-GPMRS, an extension to MR-GPSRS.

Experiments have shown that, for large data sets with large skylines, MR-GPMRS is efficient at utilizing multiple reducers to compute the global skyline. For some large skyline data sets, it is even able to return the skyline with a relatively low runtime of a few minutes, where the single reducer

solution, MR-GPSRS, can not complete in a reasonable amount of time. For data sets with low skyline percentages, the extra overhead in MR-GPMRS, however, is not worth it, meaning that the original method with one reducer in MR-GPSRS has a lower runtime.

A simple and effective solution has also been found to the tuning problem. The algorithm makes a guess of a good PPD and does the initial work of generating and pruning the bitstring multiple times for different bitstrings based on the PPD the algorithm guessed, and other PPDs close to it. After finishing the bitstrings, the algorithm determines the best one and continues with it. This solution is not especially elegant, relying on approximation and redundant work. However, the redundant work, generating and pruning the bitstring, is computationally light and it is barely noticeable on the runtime. The method is effective in that it works well enough to be able to, for each of several tested data sets, choose the optimal PPD. The method was used for all experiments. Whether the PPDs it produced for the algorithm were optimal was not tested for in all data sets. Considering how the runtime scaled, however, the produced PPDs were at least good estimates, never being bad enough for it to be noticeable.

Throughout the semester, general optimisations and refinements of the algorithms were made. Some of the implementation of MR-GPSRS, and by extension MR-GPMRS, was improved to be more efficient, and several places were identified where improvements could be made in the algorithm design. Algorithms from the literature was implemented to use for comparison, and an effort was made to implement them to be as efficient as possible.

A cost estimation for the algorithm was made. The original purpose of the cost estimation was to estimate how many dominance checks between tuples the algorithm makes. It proved to be problematic because of the high number of variables in the algorithm. A method was found to calculate the worst case number of times dominance checks between partitions are made. The most important factor for the runtime in the algorithm, the number of tuples per partition, could not be accounted for, given its unpredictability. The consequence of this is that a lower estimated number of partition dominance checks does not necessarily translate to a lower runtime. The cost estimation was compared with measured numbers and it is shown to be accurate for data sets with uniform data distribution.

1.3 Reflection

With regards to the balance of time usage, what could have been better is how much time was spent fine tuning the algorithm after it was practically finished and ready for use. Less time spent on that could have been used to make more experiments. Specifically, more testing on how many reducers were optimal to use. Throughout the experiments, a number of reducers

equal to the number of nodes in the cluster was used. Some testing with varying number of reducers was performed, which revealed that using more reducers than the number of nodes could actually provide a better runtime, the nodes being able to process more than one reducer at a time by utilising multiple cores. The experiments only showed a very small reduction in runtime when more reducers were used than there were nodes, but further investigation is warranted to find out what the optimal number of reducers is for different data sets.

Future work on the algorithm could be to further investigate the way reducer groups are generated. Currently it is done by balancing the workload of the reducers. An alternative method was found that optimised communication cost, but it was consistently showing worse runtimes than the method based on workload balancing. It would be interesting to find a more effective way to balance reducer groups. The reducer groups produced did end up being quite well balanced, only causing a difference in runtime between reducers of a few percent. This does not leave a lot of room for improvement on that area. What could be improved is the communication cost. As mentioned, an alternative method optimised communication cost, but maybe a method that compromises between load balancing and communication cost optimization can be found.

The current solution to the tuning problem is naive in the way that it assumes a uniform distribution among the partitioning when deciding on a grid configuration. In the case of a uniformly distributed data set this might be optimal, but it is not optimal for anti-correlated data sets. A future direction for research could be to develop a more intelligent way of deciding the grid configuration based on the total amount of partitions and the possibilities of pruning using the grid. It would also be interesting to further investigate how the current method performs. Manually finding the optimal PPD to a data set is very time consuming, given that multiple PPDs must be tested for. This was done for some data sets to gauge the precision of the method for choosing the PPD. It would, however, be interesting to manually find the optimal PPD for more data sets to verify that the PPD the method chooses is consistently optimal.

Chapter 2

Scientific Article

Efficient Skyline Computation for Large Volume Data in MapReduce Utilising Multiple Reducers

Kasper Mullesgaard
Aalborg University
Aalborg, Denmark

kmulle08@student.aau.dk

Jens Laurits Pedersen
Aalborg University
Aalborg, Denmark

jlpe08@student.aau.dk

ABSTRACT

A skyline query is useful for extracting a complete set of interesting tuples from a large data set according to multiple criteria. The sizes of data sets are constantly increasing and the architecture of backends are switching from single node environments to cluster oriented setups. Previous work has presented ways to run the skyline query in these setups using the MapReduce framework, but the parallel possibilities are not taken advantage of since a significant part of the query is always run serially. In this paper, we propose the novel algorithm MapReduce - Grid Partitioning Multiple Reducers Skyline (MR-GPMRS) that runs the entire query in parallel. This means that MR-GPMRS scales well for large data sets and large clusters. We demonstrate this using experiments showing that MR-GPMRS runs several times faster than the alternatives for large data sets with high skyline percentages.

1. INTRODUCTION

Skyline queries have a wide application domain ranging from e-commerce and quality based service selection, to stock trading and, generally speaking, any process involving multi-attribute decision making. Skyline query processing is computationally intensive. In order to handle this high computation cost, commodity computing can be used. Commodity computing is a paradigm where a high number of low cost computers are connected in a cluster to run demanding computations across multiple nodes. Commodity computing is deployed by several notable companies [1, 7, 8, 9, 11].

When using commodity clusters, the idea is to take advantage of the high number of nodes by processing queries parallelly. High fault tolerance is a requirement since machines can fail, and a larger cluster have a higher probability of machines faulting. MapReduce is designed specifically for high fault tolerance parallel computing.

The problem of this article is to develop a MapReduce algorithm that finds the skyline of a large volume data set

efficiently. Similar work on this subject has been done before but previous work presents solutions that uses only a single reducer to find the global skyline, failing to utilise the MapReduce framework to its full potential. To the best of our knowledge, applying multiple reducers to find the global skyline is unprecedented when finding the skyline for large data sets, and it is the focus of the work presented here.

As a solution to the problem, the MapReduce - Grid Partitioning Multiple Reducers Skyline (MR-GPMRS) algorithm is proposed. The basis of the MR-GPMRS algorithm is, as the name suggests, to partition the input data set using a grid. A bitstring is used to keep track of which partitions in the grid are non-empty, which makes it possible to make decisions based on the distribution of the entire data set. How dominance in the skyline query works combined with the grid partitioning scheme allows splitting the data into parts that can be processed independently of each other, which means that these parts can be processed by different reducers. The goal of MR-GPMRS is to minimise the query response time for data sets with high skyline percentages. To measure this, it is compared to other skyline query processing algorithm that applies the MapReduce framework.

The rest of this paper is organized as follows: Section 2 contains the preliminaries. In Section 3 we introduce MapReduce - Grid Partitioning Single Reducers Skyline (MR-GPSRS): Our bitstring based grid partitioning algorithm for skyline query processing in the MapReduce framework. In Section 4 we introduce MR-GPMRS, a novel extension to MR-GPSRS that allows processing skyline queries in MapReduce using multiple reducers. In Section 6 we present experimental evaluation of the proposed solutions compared with algorithms from the literature, and finally, in Section 7, we conclude the paper and propose future directions for research.

2. PRELIMINARIES

In this section, the skyline query and the MapReduce framework is described. A table of common symbols used throughout this paper is shown in Table 1.

2.1 The Skyline Query

Given a set of multi-dimensional tuples R , the skyline query returns a set of tuples S_R , such that S_R consists of all the tuples in R that are not dominated by any other tuple in R [3].

Definition 1. A tuple r_i dominates another tuple r_j , denoted by $r_i \prec r_j$, if and only if, for all dimensions, the value

Table 1: A list of common terms.

Symbol	Interpretation
R	A set of tuples
S_R	The skyline of the set of tuples R
t	A tuple
n	Partitions per dimension (PPD)
d	Dimensionality
p	A partition of the data
P	A set of partitions
BS	A bitstring
IG	A group of independent partitions

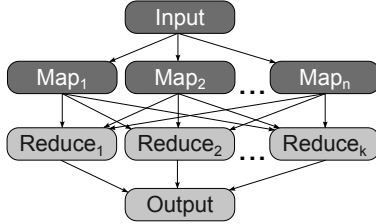


Figure 1: This figure shows the MapReduce process. The input is split between the mappers. The mappers process their input split and output the results to the reducers. The reducers process the results from the mappers, generating the final output.

of r_i is not worse than the corresponding value of r_j , and for at least one dimension, the value of r_i is better than the value of r_j [3].

Whether a value is better or worse than another value is determined by the configuration of the skyline query. Typically, a value v_1 has to be either larger or smaller than another value v_2 for v_1 to be better than v_2 . In this paper it is assumed that a smaller value is better.

2.2 The MapReduce Framework

MapReduce is a framework for distributed computing. It is based on a Map and a Reduce function [6]. The Map function is invoked for each record in the input file and it produces a list of key-value pairs. The Reduce function is then invoked once for each unique key and the associated list of values. This produces key-value pairs that are the result of the MapReduce job, i.e., $Map(k1, v1) \rightarrow list(k2, v2)$ and $Reduce(k2, list(v2)) \rightarrow list(k3, v3)$. Several MapReduce jobs can be chained together, later phases being able to refine and/or use the results from earlier phases. The MapReduce process is illustrated in Figure 1.

A distributed file system is used to store the data processed and produced by the MapReduce job. The input file(s) is split up, stored, and possibly replicated on the different nodes in the cluster. The nodes are then able to access their local splits when processing data. When the data from the Map function has been processed by the different nodes, the results are shuffled between the nodes so the required data can be accessed locally when the Reduce function is invoked.

It can be necessary to replicate some data across all nodes. In Hadoop [2], the implementation of MapReduce used for this paper, the Distributed Cache can be used for this purpose. In the beginning of a MapReduce job, data written

to the Distributed Cache is transferred to all nodes, making it accessible in the Map and Reduce functions. This paper assumes that the Distributed Cache, or something similar, is available.

2.3 Skyline Query Processing in MapReduce

In an article by Zhang et al. [12] skyline algorithms are adapted for the MapReduce framework. Three different algorithms are presented: MapReduce - Block Nested Loop (MR-BNL), MapReduce - Sort Filter Sort (MR-SFS), and MR-Bitmap. MR-BNL uses BNL and grid partitioning. The second algorithm, MR-SFS, modifies MR-BNL with presorting, but it is shown to perform worse than MR-BNL. MR-Bitmap is based on a bitmap which is used to determine dominance. It is fast computationally but requires a large amount of disk space and is only viable for data sets with few distinct values. A single reducer is used to calculate the final resulting skyline in MR-BNL and MR-SFS. MR-Bitmap does use multiple reducers, and is the only MapReduce algorithm for finding the skyline of a data set we know of to do so. However, as mentioned, it can only handle data sets with low data distinction. In [12], it was not tested on data sets with more than ten thousand distinct values, which is below the threshold for data sets used for testing in this article.

Angular partitioning is a different partitioning technique proposed by Vlachou et al. [10]. Angular partitioning is based on making partitions by dividing the data space up using angles. The idea is based on the observation that skyline tuples are located near the origin. So by dividing the data space up using angles, skyline tuples should be distributed into several partitions while non-skyline tuples should be grouped with skyline tuples that dominates them. The technique is shown to be effective but the global skyline is found using a single node. In an article by Chen et al. [4] the angular partitioning technique is adapted to MapReduce resulting in the algorithm MapReduce - Angle (MR-Angle). The results are comparable to those in [10], and the global skyline is found using a single reducer.

3. GRID PARTITIONING BASED SINGLE REDUCER SKYLINE COMPUTATION

In this section, a skyline algorithm for MapReduce is proposed. The algorithm utilises grid partitioning and bitstrings in order to prune dominance checks between tuples.

3.1 Grid Partitioning

Grid partitioning is a method of splitting up a space where each dimension is divided into n parts, referred to as the Partitions per Dimension (PPD). This gives a regular grid of n^d partitions, termed as P , where d is the dimensionality of the data set. In the context of skyline queries, the dominating relationship between the partitions $p_1, p_2, \dots, p_n \in P$ can be exploited to exclude dominance checks between tuples. Partitions have a dominating relationship with each other similar to that between tuples. The main difference is that a dominating relationship between two partitions p_i and p_j is based on their maximum corners $p_i.max$ and $p_j.max$, and minimum corners, $p_i.min$ and $p_j.min$. The maximum corner of a partition is defined as the corner of the partition that has the highest (worst) values. Similarly, the minimum corner of a partition is defined as the corner of the partition that has the lowest (best) values.

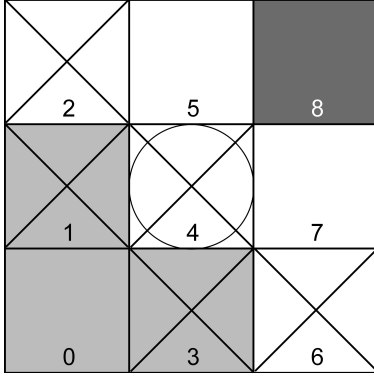


Figure 2: An example of partitions and their offsets. Non-empty partitions are marked with crosses and the dominating (dark grey) and anti-dominating (light grey) regions of the partition with offset 4 (circle) are shown.

Definition 2. A partition p_i dominates another partition p_j , denoted by $p_i \prec p_j$, if and only if $p_i.max$ dominates $p_j.min$. This ensures that all tuples in p_i dominates all tuples in p_j .

$$p_i \prec p_j \Leftrightarrow p_i.max \prec p_j.min \quad (1)$$

If this is not the case, p_i does not dominate p_j , denoted by $p_i \not\prec p_j$.

The dominating relationships between the partitions can be expressed using their individual dominating (Cui et al. [5]) and anti-dominating regions.

Definition 3. Given a partition p_i , its dominating region $p_i.DR$ contains all partitions dominated by p_i :

$$p_i.DR = \{p_j \mid p_j \in P \wedge p_i \prec p_j\} \quad (2)$$

Meanwhile, p_i 's anti-dominating region $p_i.ADR$ contains all partitions that can have tuples that dominates $p_i.max$:

$$p_i.ADR = \{p_j \mid p_j \in P \wedge p_j.min \prec p_i.max\} \quad (3)$$

Figure 2 shows an example of the dominating and anti-dominating region of the partition marked as 4 in a two dimensional data set. The non-empty partitions are marked with a cross. The dominating region of partition 4 contains partition 8 and the anti-dominating region contains partitions 0, 1, and 3.

3.2 Bitstring Representation

In grid partitioning, the only partitions of interest are those that are non-empty, i.e. $\{p_i \mid p_i \in P \wedge p_i \neq \emptyset\}$. The partitioning scheme can be represented as a bitstring $BS(0, 1, 2, \dots, n^d - 1)$ where for $0 \leq i \leq n^d - 1$:

$$BS[i] = \begin{cases} 1 & \text{if } p \neq \emptyset \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

The resulting bitstring can be constructed using either row-major order or column-major order, the only difference being how the offset of a partition in the bitstring is calculated. Column-major order is used in this paper. For example, the offset of the partitions of the two dimensional data set in Figure 2 is indicated by the digit in their lower left corner of the partitions, resulting in the bitstring 011110100.

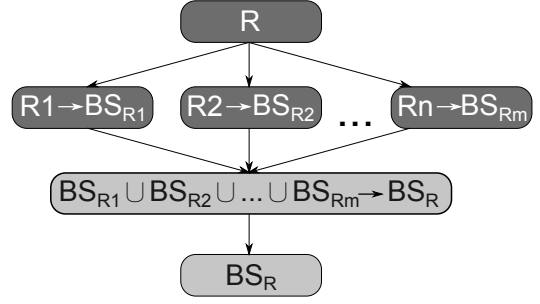


Figure 3: The data flow of the bitstring generation phase of MR-GPSRS with the data set R .

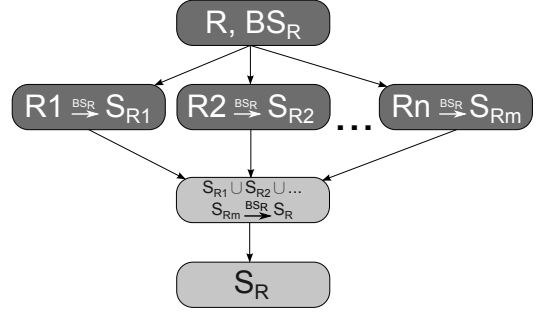


Figure 4: The data flow of the skyline computation phase of MR-GPSRS with the data set R and the bitstring BS_R .

The bitstring can be traversed to prune partitions such that fewer partitions and data tuples are involved in the skyline computation. This can be done by using the dominating relationships between partitions. If $p_i \prec p_k$ for $p_i, p_k \in P$ then the value of p_k in the bitstring is set to 0, thereby eliminating it from further consideration. If n subsets $R_1, R_2, \dots, R_n \subset R$ are partitioned with the same grid scheme, this will result in n bitstrings $BS(R_1), BS(R_2), \dots, BS(R_n)$. Two or more of these bitstrings can be merged using bitwise or, and if $R_1 \cup R_2 \cup \dots \cup R_n = R$ then $BS(R_1) \vee BS(R_2) \vee \dots \vee BS(R_n) = BS(R)$.

3.3 MR-GPSRS Algorithm

The algorithm is divided into two phases: The bitstring generation phase and the skyline computation phase.

In the mappers of the bitstring generation phase, shown in Algorithm 1, a bitstring BS_{Ri} is initialized (line 1), the status of the partitions of the tuples in Ri are set to 1 in BS_{Ri} (lines 2-5), and all mappers send their BS_{Ri} to a single reducer (line 6). In the reducer (Algorithm 2) the global bitstring BS_R is initialized (line 1). A logical OR operation is then performed on the global bitstring BS_R and each of the bitstrings received from the mappers BSS (lines 2-4). The bitstring is then traversed and pruned (lines 5-7) for dominated partitions. The data flow of the bitstring generation phase of MR-GPSRS is represented in Figure 3 where it is shown how the data set R is split into subsets $[R_1, R_m]$ that are processed by the mappers into local bitstrings $[BS_{R1}, BS_{Rm}]$. The reducer then finds the global bitstring BS_R using $[BS_{R1}, BS_{Rm}]$ and outputs BS_R .

In the skyline computation phase, shown in Algorithm 3, the mappers partition their subset Ri of the data set R (lines 1-2). By using the bitstring BS_R from the bitstring genera-

Algorithm 1 *Mapper of the bitstring generation phase*

Input: A subset R_i of the data set R , the dimensionality of the data set d , and the PPD n .

Output: A bitstring BS_{R_i} of the empty and non-empty status of the partitions in the data set RS .

- 1: Initialize a bitstring BS_{R_i} with length n^d where all bits are set to 0
 - 2: **for each** $t \in R_i$ **do**
 - 3: Decide the partition p that t belongs to
 - 4: Set the bit that represents the status of p in bitstring BS_{R_i} to 1
 - 5: **end for**
 - 6: Output($null, BS_{R_i}$)
-

Algorithm 2 *Reducer of the bitstring generation phase*

Input: A set of local bistrings BSS , the dimensionality of the data set d , and the PPD n .

Output: BS_R , the bitstring of the data set R .

- 1: Initialize a bitstring BS_R with length n^d where all bits are set to 0
 - 2: **for each** $BS_{R_i} \in BSS$ **do**
 - 3: $BS_R \leftarrow BS_R \vee BS_{R_i}[i]$
 - 4: **end for**
 - 5: **for each** partition p with status 1 in BS_R **do**
 - 6: set status of partitions in $p.DR$ to 0 in BS_R
 - 7: **end for**
 - 8: Output($null, BS_R$)
-

tion phase, the tuples belonging to partitions that have been pruned are discarded (line 3). The *InsertTuple* function (Algorithm 5) is then called on the remaining tuples (line 4), where a tuple t is inserted into a partition p and t is compared with all tuples in p , such that only the local skyline of p is maintained. The mappers, after all tuples have been partitioned, compute their local skyline (lines 7-9) by calling the *ComparePartitions* function (Algorithm 4), where the dominating relationships of the partitions, as specified in Definition 3 (Section 3.1), are used to remove any tuples dominated by other tuples local to that mapper. Each mapper then outputs their set of local partitions P to a single reducer (line 10). At this point in the algorithm, P is equal to S_{R_i} , the skyline of the mappers local subset R_i of the data set R .

In Algorithm 6, the reducer receives the local skylines, in the form of a set of partition sets LS , from the mappers and merges them (lines 1-8). Merging the partitions uses the same function for inserting tuples into partitions as in Algorithm 3, where only the local skyline is maintained in each partition. The reducer then calculates and outputs S_R , the skyline of R , by iterating through the partitions and comparing them with the partitions in their anti-dominating regions (lines 9-11).

The data flow of the skyline computation phase of MR-GPSRS is represented in Figure 4 where it is shown how the data set R is split into subsets $[R_1, R_m]$ that are processed by the mappers into local skylines $[S_{R_1}, S_{R_m}]$. The reducer then finds the global skyline S_R using $[S_{R_1}, S_{R_m}]$ and outputs S_R .

Algorithm 3 *Mapper of MR-GPSRS Skyline Computation*

Input: A subset R_i of the data set R and the bitstring BS_R .

Output: A set of local partitions S_{R_i} where each partition contains local skyline tuples.

- 1: **for each** $t \in R_i$ **do**
 - 2: decide the partition p in the set of local partitions P that t belongs to
 - 3: **if** status of p in BS_R is 1 **then**
 - 4: $p \leftarrow \text{INSERTTUPLE}(t, p)$
 - 5: **end if**
 - 6: **end for**
 - 7: **for each** $p \in P$ **do**
 - 8: $p \leftarrow \text{COMPAREPARTITIONS}(p, P)$
 - 9: **end for**
 - 10: Output($null, P$)
-

Algorithm 4 *ComparePartitions(partition p , set of partitions P)*

Input: A partition p and a set of partitions P

Output: Returns p such that all tuples in p dominated by a tuple in any partition in P are removed.

- 1: $ADR = p.ADR \cap P$
 - 2: **for each** $p' \in ADR$ **do**
 - 3: remove from p tuples that are dominated by tuples in p'
 - 4: **end for**
 - 5: **return** p
-

Algorithm 5 *InsertTuple(tuple t , partition p)*

Input: A tuple t and a partition p

Output: Returns p such that it contains t if t is not dominated by any tuples in p . If t dominates any tuples in p , they are removed.

- 1: $check = true$
 - 2: **for each** $t' \in p$ **do**
 - 3: **if** $t < t'$ **then**
 - 4: remove t' from p
 - 5: **end if**
 - 6: **if** $t' < t$ **then**
 - 7: $check = false$
 - 8: BREAK
 - 9: **end if**
 - 10: **end for**
 - 11: **if** $check$ **then**
 - 12: add t to p
 - 13: **end if**
 - 14: **return** p
-

Algorithm 6 *Reducer of MR-GPSRS Skyline Computation*

Input: The set of local skylines from all the mappers LS in the form of a set of partition sets.

Output: The global skyline S_R .

```

1: for each  $P \in LS$  do
2:   for each  $p \in P$  do
3:     for each  $t \in p$  do
4:       decide the partition  $p'$  in the set of global
         partitions  $P_G$  that  $t$  belongs to
5:        $p' \leftarrow \text{INSERTTUPLE}(t, p')$ 
6:     end for
7:   end for
8: end for
9: for each  $p' \in P_G$  do
10:   $\text{COMPAREPARTITIONS}(p', P_G)$ 
11: end for
12: for each  $p' \in P_G$  do
13:   $\text{Output}(\text{null}, p')$ 
14: end for

```

3.3.1 Choosing the Number of Partitions per Dimension

The PPD is a parameter that is significant for the performance of the algorithm. The reason the PPD is important is that it determines the number of Tuples per Partition (TPP). TPP is important since if there are too few TPP, then the process of comparing each set of partitions is not worthwhile compared to checking the tuples in the partitions. Conversely, if there are too many TPP, the grid is too rough and the number of partitions that can be pruned when comparing partitions becomes less than optimal.

What the optimal PPD is depends on the cardinality, distribution, and dimensionality of the data set, as well as the number of active mappers. Without having extensive knowledge of the algorithm and the data set, choosing the optimal PPD, or even a good PPD, is guesswork. To avoid this, an extension to MR-GPSRS is proposed where the algorithm chooses the PPD itself.

The extension is based on a guess of a good PPD n made in the mappers. This guess is based on the data sets cardinality c , dimensionality d , and the desired TPP. The number of TPP in a given grid can be approximated with the following expression:

$$\frac{c}{n^d} = \text{TPP} \quad (5)$$

From this, n can be isolated:

$$\sqrt[d]{\frac{c}{\text{TPP}}} = n \quad (6)$$

It is then necessary to determine the desired TPP. For the data sets in this article, $100/d$ was found to be a good number. This takes into account the number of dimensions, which affects the time it takes to compare them. What this number should be in different setups might vary. From the guess n , several bitstrings are generated. For example, a bitstring based on n and then four more bitstrings based on PPD values 1 and 2 higher and lower than n . The values should only be used if $2 \leq n \wedge n^d < c$.

The set of bitstrings generated by each mapper, as well as the number of points each mapper processed, is then send

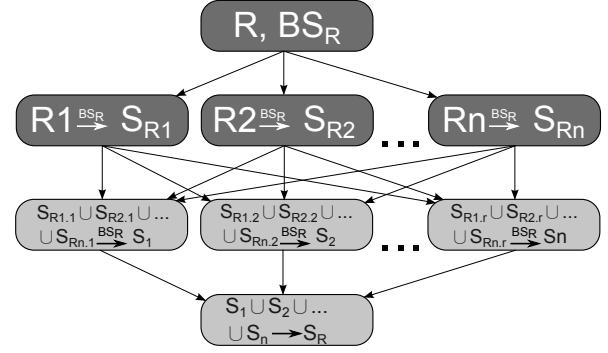


Figure 5: The data flow of the second phase of MR-GPMRS with the bitstring BS_r of the data set R

to a reducer. To get the global bitstrings, the reducer runs a logical OR operation on the bitstrings from the different mappers that are based on the same n . The reducer then makes an estimate on the number of TPP for each bitstring by dividing the number of non-empty partitions, taken from the bitstring, with the number of tuples processed by each mapper. The reducer can then estimate the remaining tuples each bitstring would have after pruning by using the estimated TPP and the difference in the number of set bits in the bitstrings before and after pruning. The remaining number of tuples and the number of non-empty partitions in each bitstring after pruning are then used by the reducer to make a final estimate of the TPP in each bitstring after pruning. The one that is closest to the desired TPP is then chosen as the final global bitstring used in the rest of the algorithm. This extension is used for the experiments in Section 6.

4. GRID PARTITIONING BASED MULTIPLE REDUCERS SKYLINE COMPUTATION

In this section, an extension to MR-GPSRS is proposed that utilises multiple reducers. MR-GPSRS relies on a single reducer for computing the global skyline, which increasingly becomes a bottleneck when the skyline of the data set becomes larger.

This bottleneck is alleviated by utilising the fact that the grid partitioning technique can be used to identify subsets of partitions for which the skyline can be computed independently, allowing the use of multiple reducers.

4.1 Skyline Query Processing Using Grid Partitioning With Multiple Reducers

In the current methods of computing skyline in MapReduce, the final step of the algorithms require the local skylines from the mappers to be merged into the global skyline by a single reducer. This is due to the inability of mappers to communicate with each other, thereby not having a global awareness of the data set, and that with some partitioning methods, it is not possible to distribute partitions into groups that can be processed independently.

The issue of communication between the mappers is addressed by utilising the bitstring to ensure that the necessary information is available to the mappers. The issue of identifying independent groups of partitions is addressed by

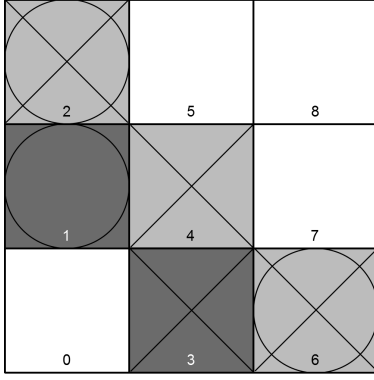


Figure 6: An example showing the distribution of non-empty partitions on two mappers (cross and circle), the partitions belonging to an independent group (grey), and the replicated partitions (dark grey).

utilising the anti-dominating relationship between the partitions in the grid partitioning scheme. The combination of these methods allows the mappers in the algorithm to unanimously decide how the partitions are to be sent to multiple reducers. In order for the mappers to output their partitions to reducers they need to be aware of which partitions are non-empty. A single mapper $mapper_1$ is aware of which partitions are non-empty in its subset. This is, however, not enough for the mapper to decide how the partitions are to be sent, as the distribution of non-empty partitions in another mapper $mapper_2$ might be different. The decision on how to send the partitions made by $mapper_1$ might be different than that made by $mapper_2$, which results in a partition p not being compared with all the partitions in its global anti-dominating region $p.ADR$.

For example, Figure 6 illustrates the non-empty partitions of the mappers $mapper_1 = \{p_1, p_2, p_6\}$ and $mapper_2 = \{p_2, p_3, p_4, p_6\}$. A decision is made by $mapper_1$ to send the partitions p_1 and p_2 to $reducer_1$ and partition p_6 to $reducer_3$. Meanwhile, $mapper_2$ decides to send partition p_2 to $reducer_1$, partitions p_3 and p_4 to $reducer_2$, and partitions p_3 and p_6 to $reducer_3$. In this example, $reducer_2$ would need partition p_1 because $p_1 \in p_4.ADR$. Since $mapper_1$ does not know that p_4 is non-empty, however, it has no way of knowing that it should send p_1 to $reducer_2$. A global bitstring is a way to resolve this problem, since it allows the mappers to know the empty or non-empty status of all partitions.

4.1.1 Grouping Partitions

The grid partitioning scheme allows for identifying *independent groups*: Sets of partitions that can be processed independently to obtain the skyline of those partitions.

Definition 4. A set of partitions IG is independent if and only if the following holds:

$$\{\forall p \in IG \mid p.ADR \subseteq IG\} \quad (7)$$

Independent groups makes it possible for the combined output of multiple reducers to be the global skyline while only having each reducer process a subset of the data set. The following lemma provides a general way to identifying *independent groups*.

Lemma 1. Any group of partitions that consists of a partition p , and the partitions in $p.ADR$, is independent.

PROOF. Consider three partitions p_1, p_2, p_3 that are part of the same regular grid for which the following statements hold:

$$p_2.min \prec p_1.max \quad (8)$$

$$p_3.min \prec p_2.max \quad (9)$$

Since the grid partitions p_1, p_2, p_3 belong to is regular, from Statements 8 and 9, it follows that:

$$p_3.min \prec p_1.max \quad (10)$$

Considering Definition 3 (Section 3.1), from Statements 8, 9, and 10 it follows that:

$$p_2 \in p_1.ADR \wedge p_3 \in p_2.ADR \Rightarrow p_3 \in p_1.ADR \quad (11)$$

which can be generalized as:

$$p_2 \in p_1.ADR \Rightarrow p_2.ADR \subseteq p_1.ADR \quad (12)$$

Considering Definition 4, it follows that the set of partitions P is independent if:

$$P = p_1.ADR \cup p_1 \quad (13)$$

□

Generating *independent groups* should not be done using arbitrarily chosen partitions and their anti-dominating regions, since this does not exclude the possibility of the *independent groups* being subsets of each other. One way to identify *independent groups* from a set of partitions P , that cannot be subsets of each other, is by using the *maximum partitions* in P .

Definition 5. A partition $p_{max} \in P$ is a *maximum partition* if and only if the following holds:

$$\{\forall p \in P \mid p_{max} \notin p.ADR\} \quad (14)$$

When a group of partitions consists of a *maximum partition* p_{max} , and the partitions in $p_{max}.ADR$, it is independent and it cannot be a subset of another *independent group*.

For example, consider Figure 6. The partition p_2 is a *maximum partition* because it is not in the anti-dominating region of another partition. Making an *independent group* IG_1 from p_2 and $p_2.ADR$ gives $IG_1 = \{p_1, p_2\}$. Similarly, partitions p_4 and p_6 are also *maximum partitions* and making *independent groups* IG_2, IG_3 from p_4 and $p_4.ADR$ and p_6 and $p_6.ADR$ respectively, gives $IG_2 = \{p_1, p_3, p_4\}$ and $IG_3 = \{p_3, p_6\}$. Since these groups are made based on *maximum partitions*, they are not subsets of each other.

It is necessary to replicate some partitions among the *independent groups* as they lie in the anti-dominating regions of partitions in multiple groups, e.g. p_1, p_3 . The skyline tuples in a replicated partition are only output by one of the reducers to which the partition is sent. The mappers decides which reducers are responsible for outputting the replicated partitions.

4.1.2 Merging Independent Groups

The number of *maximum partitions* in a data set can be high, and therefore the number of *independent groups* that can be generated will also be high. If the number of *independent groups* is higher than the number of cluster nodes, multiple *independent groups* will be sent to the same nodes,

and because partitions can be present in multiple *independent groups*, this will cause partitions to be sent to the same nodes multiple times.

Consider the previous example from Figure 6 with the three *independent groups* $IG_1 = \{p_1, p_2\}$, $IG_2 = \{p_1, p_3, p_4\}$, $IG_3 = \{p_3, p_6\}$. In a scenario where IG_1 and IG_2 are sent to the same node, for example, partition p_1 is sent twice. It is notable that the overlap between partitions becomes more prominent in higher dimensions where the number of replicated partitions increases, since the dimensionality of the anti-dominating regions of the partitions also increases.

The *independent groups* can be merged, however, to avoid sending the same partitions to the same nodes multiple times. The two groups IG_1 and IG_2 can be merged to form the group $IG_{merged} = \{p_1, p_2, p_3, p_4\}$. These merged groups are then able to be sent to reducers without any duplication in their list of partitions, and are called *reducer groups*. The merging of *independent groups* influences both the communication cost and the balancing of computation cost between the reducers. One method of merging is based on optimizing the communication cost: *Independent groups* that have the most partitions in common are merged. This method, however, does not guarrenty any balance of the computations among the reducers as this could leave some *reducer groups* with more unique partitions than other groups. Consider the previous example, the *reducer group* IG_{merged} consists of 3 unique partitions and IG_3 consists of 1 unique partition. This would result in a larger amount of computations for IG_{merged} , assuming uniform amount of tuples in the partitions, as it is forced to compute the skyline for the 3 partitions as they are not present in any other *reducer group*. Preliminary tests have shown that a merging method based on balancing the computations cost between the reducers performed the best, and it was therefore the one used throughout this paper.

How many *reducer groups* that should be generated depends on the size of the data set, the size of the cluster, and the available memory in the nodes. The reducer part of the skyline generation phase requires that the skyline of the local data set is kept in memory.

This means that if the number of *reducer groups* is set too low compared to the size of the data set, the memory may overflow which will cause the runtime to increase significantly. If the number of *reducer groups* is set too high, the same partitions can be sent to the same node several times, or be sent to multiple nodes unnecessarily, causing additional communication. A balance between the two is desirable, such that memory does not overflow and communication cost is not unnecessarily high. In this paper, the number of generated *reducer groups* is set to be equal to the number of nodes in the cluster r . This is fitting since the local skylines of the data sets tested for fits in the memory of the nodes and any unnecessary communication cost is avoided. For larger clusters or larger data sets, generating a number of *reducer groups* equal to the number of nodes might not be optimal.

4.1.3 Outputting Replicated Partitions

Since partitions are replicated in different *reducer groups*, it is necessary to control which reducers output the replicated partitions. When a reducer is responsible for outputting the skyline tuples in a partition, the *reducer group*

of that reducer is referred to as being responsible for that partition. A *reducer group* is responsible for all the partitions unique for that *reducer group*. If a partition is present in more than one *reducer group*, one of the *reducer groups* in which it is present is designated as being responsible for that partition. The *reducer group* chosen to be responsible for a replicated partition is based on a calculation of how many comparisons are made between the partitions it contains. The number of comparisons necessary for each partition the *reducer group* is responsible for is calculated based on the bitstring and added up. In order to balance the computation cost for each *reducer group*, the *reducer group* with the lowest number of calculated partition comparisons is made responsible for a replicated partition until responsibility of every partition has been given to some *reducer group*. Since the mappers use the bitstring to form *reducer groups*, the allocation of responsibility is identical across all mappers.

4.2 GPMRS Algorithm

The MR-GPMRS algorithm has two phases; the bitstring generation phase and the skyline computation phase, where the bitstring generation phase is the same as in MR-GPSRS. The skyline computation phase of MR-GPMRS, Algorithm 7, is the same as the skyline computation phase of MR-GPSRS until line 9. There after the *independent groups* are found (line 10), grouped together (line 11), and responsibility of replicated partitions are assigned (lines 12-16). The *reducer groups* are then sent to the reducers (lines 17-21).

Algorithm 7 Mapper of MR-GPMRS

Skyline computation

Input: A subset R_i of the data set R , the bitstring BS_R , and the number of reducers r .

Output: A set of *reducer groups* RG containing local partitions with the local skyline.

```

1: for each  $t \in R_i$  do
2:   Decide the partition  $p$  in the set of local partitions
    $P$  that  $t$  belongs to
3:   if status of  $p$  in  $BS_R$  is 1 then
4:      $p \leftarrow \text{INSERTTUPLE}(r, p)$ 
5:   end if
6: end for
7: for each  $p \in P$  do
8:    $p \leftarrow \text{COMPAREPARTITIONS}(p, P)$ 
9: end for
10:  $IG \leftarrow \text{INDEPENDENTGROUPS}(P, BS_R)$ 
11:  $RG \leftarrow \text{REDUCERGROUPS}(IG, r)$ 
12: while  $BS_R$  contains set bits do
13:    $rg_{min} \leftarrow$  a reference to  $rg \in RG$  with
   the lowest amount of computations
14:   Assign responsibility of a single  $p \in rg_{min}$ 
   to  $rg_{min}$ 
15:   Set index of  $p$  in  $BS_R$  to 0
16: end while
17:  $i = 0$ 
18: for each  $rg \in S_{R_i} \leftarrow RG$  do
19:   Output( $i, rg$ )
20:    $i++$ 
21: end for

```

In order to maintain consistency throughout the mappers, the bitstring is used to generate the *independent groups*. In Algorithm 8, the independent groups are generated by

traversing the bitstring BS_R in reverse in order to find the *maximum partitions* (line 2). For each maximum partition p_{max} an *independent group* is created consisting of p_{max} and the partitions in $p_{max}.ADR$ (lines 3-4). The *independent group* ig is then added to the set of *independent groups* IG (line 5) and the bitstring indexes of the partitions in the *independent group* ig are cleared from the bitstring BS_R (lines 6-7), so that the next traversed bit will be a new maximum partition. The result is a set of independent groups IG (line 9). Consider the previous example of Figure 6 with the independent groups $IG_1 = p_1, p_2$, $IG_2 = p_1, p_3, p_4$ and $IG_3 = p_3, p_5$, the result of the algorithm is the set of independent groups $IG = IG_1, IG_2, IG_3$.

Algorithm 8 *IndependentGroups(P, BS_R)*

Input: A set of partitions P and a bitstring BS_R .

Output: A set of independent groups IG .

```

1: while  $BS_R$  contains set bits do
2:    $p_{max} \leftarrow$  the  $p \in P$  represented by the last set bit in  $BS_R$ 
3:    $ig \leftarrow p_{max}$ 
4:   add  $p_{max}.ADR$  to  $ig$ 
5:   add  $ig$  to  $IG$ 
6:    $BS_R \leftarrow BS_R \wedge \neg$  bitstring of  $p_{max}$ 
7:    $BS_R \leftarrow BS_R \wedge \neg$  bitstring of  $p_{max}.ADR$ 
8: end while
9: return  $IG$ 
```

Algorithm 9 *ReducerGroups(IG, r)*

Input: A set of independent groups IG and the required number of *reducer groups* r .

Output: A set of *reducer groups* RG .

```

1:  $RG \leftarrow r$  number of empty reducer groups
2: while  $IG \neq \emptyset$  do
3:    $rg_{min} \leftarrow$  a reference to the  $rg \in RG$  with the lowest amount of unique partitions
4:   move the partitions in the largest  $ig_{max} \in IG$  to  $rg_{min}$ , ignoring duplicate partitions
5: end while
6: return  $RG$ 
```

In Algorithm 7 the *independent groups* IG are grouped together into r *reducer groups*, where r is the number of reducers (line 11). In Algorithm 9, the r *reducer groups* are constructed by continually moving *independent groups* from IG to the *reducer group* rg_{min} with the lowest amount of unique partitions (lines 2-5). In the start r empty *reducer groups* are initialized as the *reducer groups* (line 1). Then rg_{min} is found and the largest (in number of partitions) *independent group* ig_{max} is moved from IG to rg_{min} (lines 3-4). The procedure then returns the set of *reducer groups* RG (line 6). Continuing the previous example, with the set of *independent groups* $IG = \{IG_1, IG_2, IG_3\}$. In this example, 2 *reducer group* are generated, so 2 empty *reducer groups* rg_1 and rg_2 are initialized. The *reducer group* with the lowest number of partitions is chosen, which can be either one, so the first RG_1 is chosen. The *independent group* with the highest number of partitions, IG_2 , is then added to the *reducer group*, $rg_1 = \{p_1, p_3, p_4\}$. Then the *reducer group* with the lowest number of partitions is chosen again, rg_2 , and the *independent group* with the most partitions is added

to it. This can be either IG_1 or IG_3 since they contain the same number of partitions, so the first one is chosen, $rg_2 = \{p_1, p_2\}$. The *reducer group* with the lowest number of partitions is still RG_2 so it is chosen again and the last *independent group* is added to it, $rg_2 = \{p_1, p_2, p_3, p_6\}$. The set of *independent groups* IG is empty so the *reducer groups* are finished being generated, $RG = \{rg_1, rg_2\}$.

Responsibility of tuples are assigned to *reducer groups* in Algorithm 7 (lines 12-16). This is done by calculating the increase in number of required dominance checks between partitions the partition would cause if added to the *reducer groups*, which is how many more times line 3 in Algorithm 4 would be run. The partition is then assigned to the *reducer group* with the lowest number of added partition comparisons, rg_{min} .

Algorithm 10 *Reducer of MR-GPMRS*

Skyline computation

Input: A set of *reducer groups* LS containing subsets of the local skylines send to the same reducer from different mappers.

Output: The partitions the *reducer groups* in LS are responsible for.

```

1: for each  $RG \in LS$  do
2:   for each  $p \in RG$  do
3:     for each  $t \in p$  do
4:       decide the partition  $p'$  in the set of global partitions  $P_G$  that  $t$  belongs to
5:        $p' \leftarrow \text{INSERTTUPLE}(t, p')$ 
6:     end for
7:   end for
8: end for
9: for each  $p \in P_g$  do
10:  if  $p$  is responsible then
11:     $p \leftarrow \text{COMPAREPARTITIONS}(p, P_G)$ 
12:  end if
13: end for
14: for each responsible  $p \in P_G$  do
15:   Output( $null, p$ )
16: end for
```

Consider the previous example $RG = \{rg_1, rg_2\}$. *Reducer group* rg_1 starts with taking responsibility of the partition it contains that adds the lowest amount of partition comparisons. The number of comparisons a partition requires is equal to the number of partitions in its anti-dominating region. So the partitions in rg_1 with the lowest amount of comparisons are p_1 and p_3 , and it chooses the first one p_1 . Partition p_1 requires no comparisons so it chooses again and takes responsibility for p_3 , which also does not require any comparisons. So rg_1 chooses again and takes responsibility for partition p_4 that requires 2 comparisons. *Reducer group* rg_1 has no more partitions to take responsibility for, so rg_2 takes responsibility of its partitions in order p_2, p_6 and ends up requiring 2 comparisons.

In Algorithm 10, each reducer receives a set of *reducer groups* LS , containing their subset of the local skyline from one or more mappers, which are then merged into a single group MRG (lines 1-8). The reducers then calculates their subset of the skyline of R , by iterating through the partitions they are responsible for and comparing them with the partitions in their anti-dominating regions (lines 9-11). The result is then output (lines 14-16). The data flow of the skyline

Table 2: A list of common terms.

Symbol	Interpretation
$p_{total}(n, d)$	Partitions in a grid
$p_{rem}(n, d)$	Remaining partitions in a grid after pruning dominated partitions
$p_{dom}(n, d)$	Partition dominance checks for a single partition
$s(n, d)$	Partition dominance checks for a single surface in a grid
$g_{mapper}(n, d)$	Partition dominance checks for a single mapper
$g_{reducer}(n, d)$	Partition dominance checks for the reducer with the most partition dominance checks

computation phase of MR-GPMRS is represented in Figure 5 where it is shown how the data set R is split into subsets R_1, R_2, \dots, R_n that are processed by the mappers into local skylines $S_{R_1}, S_{R_2}, \dots, S_{R_n}$. The reducers then find subsets of the global skyline S_1, S_2, \dots, S_n using subsets of the local skylines $\{\{S_{R_1.1}, \dots, S_{R_1.n}\}, \{S_{R_2.1}, \dots, S_{R_2.n}\}, \dots, \{S_{R_n.1}, \dots, S_{R_n.n}\}\}$ which are then combined to produce the global skyline S_R as output.

5. COST ESTIMATION

In this section estimations of the cost of MR-GPMRS are made. A table of common terms is shown in Table 2.

5.1 Partition-wise Dominance Tests Estimate

The purpose of this section is to estimate the number of dominance checks performed between partitions in the MR-GPMRS algorithm. Specifically, what is estimated is how many times the line 3 in Algorithm 4 is executed. Due to several uncertainties of the algorithm, it is necessary to make some assumptions.

- **Every partition in every mapper is non-empty.** The data distribution is an important factor of the estimation. It is necessary to assume that there are no empty partitions in order to predict the required number of comparisons. It is comparable to uniform data distribution where it can be expected that most partitions are non-empty.
- **The dominance checks between partitions in the mappers do not lower the number of non-empty partitions.** In practice, the partition checks in the mappers are likely to leave some partitions empty. This is unpredictable and is therefore not accounted for.

These assumptions mean that the scenario for which the estimations are made is a worst case scenario for a data set with a uniform data distribution. When every partition is non-empty, after the grid has been pruned, the location and amount of the remaining partitions is predictable using the dimensionality d and the PPD n . A d dimensional grid has a number of $d - 1$ dimensional surfaces equal to $d \times 2$. Half of these surfaces, i.e. d surfaces, are filled with remaining partitions. The remaining part of the other half of the surfaces, as well as the rest of the partitions, are

dominated. For example, consider Figure 6. In this 2 dimensional, 3 PPD grid, there are $d \times 2 = 4$ surfaces with a dimensionality of $d - 1 = 1$. These four surfaces consists of the partitions $surf_1 = \{p_2, p_1, p_0\}$, $surf_2 = \{p_0, p_3, p_6\}$, $surf_3 = \{p_6, p_7, p_8\}$, and $surf_4 = \{p_8, p_5, p_2\}$. If every partition were non-empty, the partitions p_4 , p_5 , p_7 , and p_8 would be dominated and pruned. This would leave $d = 2$ intact surfaces, $surf_1$ and $surf_2$. There is an overlap between the surfaces that must be considered. In this case, the overlap between the remaining surfaces $surf_1$ and $surf_2$ is p_0 .

The number of remaining partitions after pruning a grid where every partition is non empty $p_{rem}(n, d)$ can be calculated by finding the total number of partitions in a grid $part_{total}(n, d)$ and subtracting a grid one PPD smaller:

$$p_{total}(n, d) = n^d \quad (15)$$

$$p_{rem}(n, d) = part_{total}(n, d) - part_{total}(n - 1, d) \quad (16)$$

From the previous example, the pruned partitions, p_4 , p_5 , p_7 , and p_8 , can be contained by a $d = 2$ dimensional $n - 1 = 2$ PPD grid. This means the the number of remaining partitions after pruning a 2 dimensional, 3 PPD grid can be calculated as $3^2 - 2^2 = 5$. The dominance checks to be done by a single partition p depends on its anti-dominating region. A partition p_i performs dominance checks against another partition p_j if $p_j \in p_i.ADR$. The number of dominance checks $p_{dom}(n, d)$ for a partition p is equal to its grid position values multiplied with each other minus one:

$$p_{dom}(n, d) = p.pos.d_1 \times p.pos.d_2 \times \dots \times p.pos.d_d - 1 \quad (17)$$

The position of a partition in a grid is how many partitions, including itself, from the origin a partition is located in the different dimensions.

For example, in Figure 6 the partition p_2 , for the first dimension, has the position $p_2.pos.d_1 = 3$, and for the second dimension it has the position $p_2.pos.d_2 = 1$. Summing this up to get the number of partition checks for every partition in a surface $s(n, d)$ yields the following expression:

$$s(n, d) = \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_d=1}^n (i_1 \times i_2 \times \dots \times i_d - 1) \quad (18)$$

To get the number of partition checks for all surfaces, the overlap between surfaces, where they meet on the axes, has to be considered. The first surface is calculated as above. The second surface is also calculated as above but with the overlap between the first and the second surface removed. The overlap between the first and the third surface and the overlap between the second and the third surface has to be subtracted from the third surface, and so on. To account for this, the start index i of one of the summations is incremented for each surface that is calculated. So the number of dominance checks $g_{map}(n, d)$ between partitions in a single mapper for all surfaces of a grid is as follows:

$$s_1(n, d) = \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_{d-1}=1}^n (i_1 \times i_2 \times \dots \times i_{d-1} - 1) \quad (19)$$

$$s_2(n, d) = \sum_{i_1=2}^n \sum_{i_2=1}^n \dots \sum_{i_{d-1}=1}^n (i_1 \times i_2 \times \dots \times i_{d-1} - 1) \quad (20)$$

$$s_3(n, d) = \sum_{i_1=2}^n \sum_{i_2=2}^n \dots \sum_{i_{d-1}=1}^n (i_1 \times i_2 \times \dots \times i_{d-1} - 1) \quad (21)$$

...

$$s_d(n, d) = \sum_{i_1=2}^n \sum_{i_2=2}^n \dots \sum_{i_{d-1}=2}^n (i_1 \times i_2 \times \dots \times i_{d-1} - 1) \quad (22)$$

$$g_{map}(n, d) = \sum_{i=1}^d s_i(n, d) \quad (23)$$

For the reducer, only a single surface has to be considered. The reason for this is that each surface is an *independent group* that can be calculated individually by the reducers. The reducer with the most dominance checks is the one that has the biggest surface, which is the one where no overlap is considered. This allows the surface calculation from before to be reused when calculating the number of partition dominance checks $g_{reducer}(n, d)$ in the reducer with the most dominance checks:

$$g_{reducer}(n, d) = s_1(n, d) \quad (24)$$

6. EXPERIMENTS

In this section, the results from experimental runs of the algorithms are presented. A cluster of thirteen commodity machines have been used for the experiments. Twelve of the machines have an Intel Pentium D 2.8 GHz Core2 processor. Three of these have a single gigabyte of RAM, four of them have two, and five of them have three. The last machine has an Intel Pentium D 2.13 GHz Core2 processor and two gigabytes of RAM. The machines are connected with a 100 Mbit/s LAN connection. The operating system used is Ubuntu 12.04 and the version of Hadoop is 1.1.0. The algorithms are implemented in Java. Tests are performed on the algorithms MR-GPMRS, MR-GPSRS, MR-BNL from [12], and MR-Angle from [4]. The tests are performed on several different data sets, with varying cardinality, dimensionality, and data distribution.

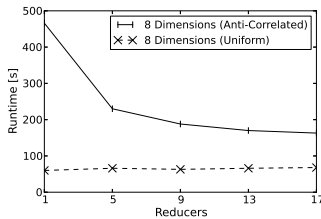


Figure 7: The graph shows the runtime of the algorithm MR-GPMRS run on an anti-correlated data set with a dimensionality of 8 and with a cardinality of 1×10^6 . The number of reducers used is varied. The result for one reducer is the runtime of MR-GPSRS.

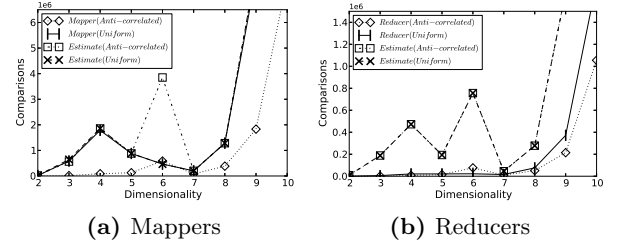


Figure 8: The estimated and actual number of dominance checks between partitions for data sets when run on a data set with a cardinality of 1×10^6 and varying dimensionality.

The cardinalities are 1×10^5 , 5×10^5 , 1×10^6 , 2×10^6 , and 3×10^6 . The dimensionality range is $[2 \dots 10]$. The distributions are anti-correlated and uniform. For the test runs of MR-GPMRS for comparison with the other algorithms, MR-GPMRS is set to use one reducer per node, i.e. thirteen. To see the effect of the number of reducers, MR-GPMRS has been run on select data sets with varying number of reducers.

6.1 Effect of Dimensionality

The runtime results for uniform data sets with varying dimensionality are shown in Figure 9. The results for anti-correlated data sets are shown in Figure 10.

As can be seen from the uniform data set results (Figure 9), the MR-GPMRS algorithm runs a little slower for the lower dimensions (2 – 5). For the rest of the dimensions (6 – 10), however, the runtime of MR-GPMRS increases linearly, while the runtime of the other algorithms increases exponentially. The exception is MR-GPSRS that runs faster than MR-GPMRS even for higher dimensions. The difference is marginal but consistent. The reason for these results are that for the low dimensionalities, the skyline percentage is so low that the runtime of the algorithms is dominated by the communication cost and the time it takes to read the data set. This overhead in MR-GPMRS is slightly larger than the other algorithms, which is why the runtime is a bit higher. For the higher dimensions, when the skyline percentage becomes more significant, MR-GPMRS is able to obtain the skyline more efficiently than MR-Angle and MR-BNL. MR-GPSRS is the fastest algorithm for all uniform data sets, which shows that the pruning power of the grid partitioning scheme used in both MR-GPSRS and MR-GPMRS is superior, and that using multiple reducers is not worth the extra communication cost when the skyline percentage is low.

For the anti-correlated data sets (Figure 10), MR-GPMRS is superior for all dimensions, except to MR-GPSRS which is better for the dimensionalities less than 5. It is clear that MR-GPMRS scales well for higher dimensionalities even for large cardinalities, having a runtime of less than ten minutes for a data set with a cardinality of 2×10^6 and 10 dimensions. Results for MR-Angle and MR-BNL are not shown for the higher dimensionalities since they are not able to terminate in a reasonable amount of time. MR-GPSRS is running significantly slower than MR-GPMRS for the low cardinality and high dimensionality data sets. For the higher cardinality, the difference is more clear and MR-GPSRS does not

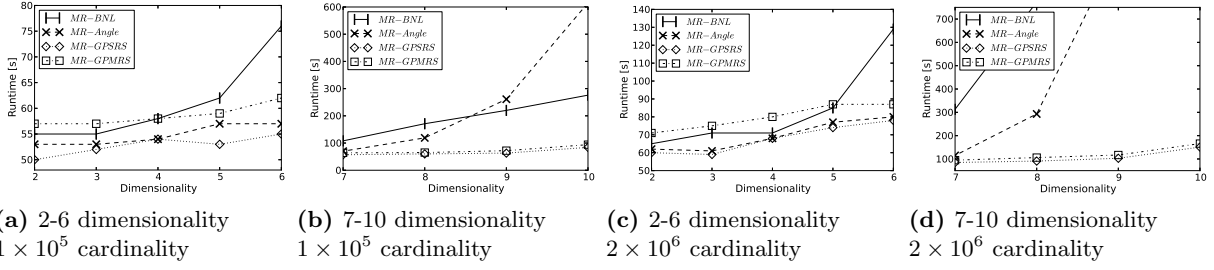


Figure 9: Each graph shows the run time of the algorithms MR-BNL, MR-Angle, MR-GPSRS, and MR-GPMRS run on uniform data sets with varying dimensionality.

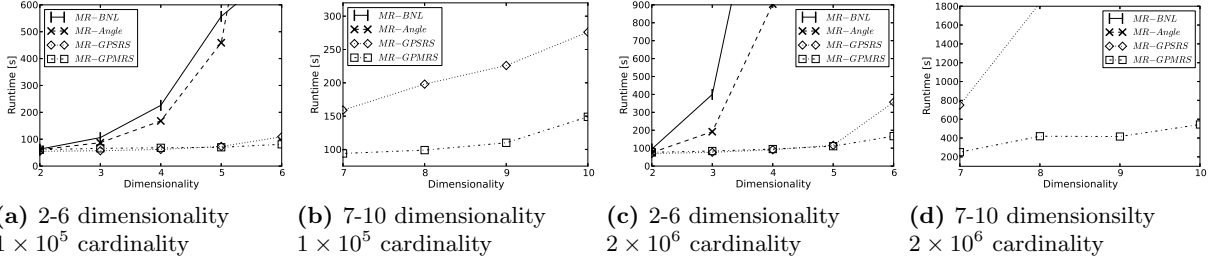


Figure 10: Each graph shows the run time of the algorithms MR-BNL, MR-Angle, MR-GPSRS, and MR-GPMRS run on anti-correlated data sets with varying dimensionality.

terminate in a reasonable amount of time for the highest dimensionalities. The reason here is that in anti-correlated data sets, the skyline is significant even when the dimensionality is low, and it increases quickly for higher dimensionalities.

This means that the ability of MR-GPMRS to find the skyline more efficiently outweighs its increased overhead. For higher dimensionalities, this is even more evident. That MR-GPSRS performs better than MR-GPMRS for the lower dimensionalities continues to show that using multiple reducers is not worth the extra communication cost when the skyline percentage is below a certain point.

6.2 Effect of Cardinality

The results for varying cardinality is shown in Figure 11. For uniform data, MR-GPMRS is slowest for all cardinalities when the dimensionality is 3, almost 20 seconds at a cardinality of 3×10^6 . MR-GPSRS has the best runtime for all cardinalities, MR-Angle tying with it for the cardinalities 1×10^6 and 3×10^6 . For the higher dimensionality of 8, MR-GPMRS and MR-GPSRS has the fastest runtimes, MR-GPSRS being a bit faster given the small skyline percentage of the uniform data where the multiple reducers are not an improvement given the extra communication cost. What these results show is that using multiple reducers consistently causes a slower runtime when the skyline percentage of the data set is low.

For anti-correlated data, MR-GPMRS and MR-GPSRS is superior for all cardinalities and both dimensionalities. For the lower dimensionality of 3, MR-GPSRS is marginally better than MR-GPMRS, but for the higher dimensionality of 8, MR-GPSRS fails to terminate in a reasonable amount of time for the highest cardinality and is consistently worse than MR-GPMRS.

6.3 Effect of Number of Reducers

Figure 7 shows how the runtime of MR-GPMRS changes when the number of reducers used is varied. What can be seen from the runtimes of the anti-correlated data set is that it lowers when the number of reducers is increased. For this particular data set the runtime is best when the number of reducers is at the highest. It is better even when the number of reducers is higher than the number of nodes. The reason for this is that Hadoop is able to utilise the multiple cores in the nodes to parallelize multiple reducers on the same node. This is not necessarily the case for other data sets, as evident from the results where MR-GPSRS outperforms MR-GPMRS. The skyline percentage needs to be high enough for the extra communication cost the multiple reducers causes for the high number of reducers to be useful. The runtime for the data set with uniform distribution almost does not change when the number of reducers is increased. There is a small increase, caused by the additional overhead of multiple reducers. This uniform data set has a much smaller skyline than the anti-correlated one, which means that the use of multiple reducers is not worthwhile given the extra overhead.

6.4 Evaluation of Estimate

Figure 8 shows the estimated number of dominance checks between partitions, as described in Section 5, compared with the measured number of partition comparisons in MR-GPMRS. The measured numbers has been taken from the mapper and the reducer that had the highest number of comparisons. The results for the mappers show that the estimate does not deviate far from the measured number when the data distribution is uniform, and in some cases there is no deviation. The estimate, however, is inaccurate

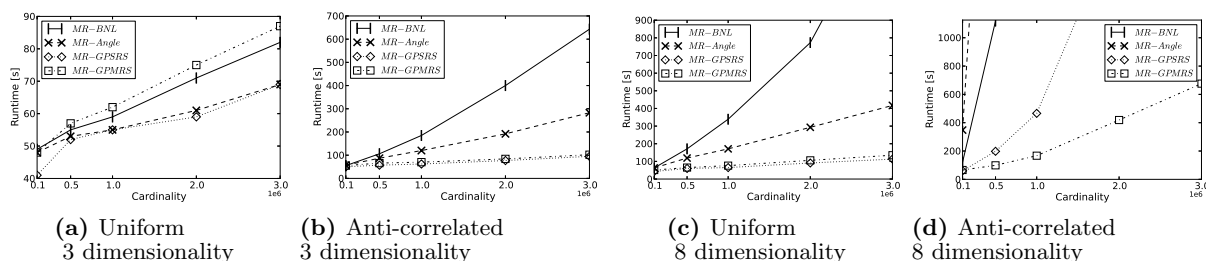


Figure 11: The graphs show the runtime of the algorithms MR-BNL, MR-Angle, MR-GPSRS, and MR-GPMRS run with a dimensionality of 3 and with varying cardinality.

when estimating the number of partition dominance checks for mappers run on the anti-correlated data set. This is to be expected as the estimation was done on the premise that the data set is uniform. It is notable that the estimated number of dominance checks is above the measured one in every case, which means that the estimate can be still be used as a worst case estimate for anti-correlated data sets. The results for the reducers show that the estimate is inaccurate for the uniform data set, is that no way was found to mathematically predict the way the *reducer groups* are generated, so instead a number was used that was ensured to be the worst case.

7. CONCLUSION

In this paper, two novel algorithms, MR-GPSRS and MR-GPMRS, for skyline query processing in MapReduce are proposed. The main feature of the algorithms are that they allow decision making across mappers and reducers. This is accomplished by using a bitstring describing the partitions empty and non-empty state across the entire data set. In addition, the common bottleneck of having the final skyline computed at a single node is avoided in the MR-GPMRS algorithm by utilizing the bitstring to partition the final skyline computation among multiple reducers.

The experiments conducted show that the algorithms proposed in this paper consistently outperforms existing algorithms and they scale well with the skyline and data set size. Which of the two proposed algorithms performs better depends on the data set. When the skyline percentage is high, MR-GPMRS performs significantly better while MR-GPSRS performs marginally better when the skyline percentage is low. Running test on a large cluster would be interesting in order investigate how well MR-GPMRS scales for large numbers of reducers.

The increased communication cost incurred by having reducers receive replicated partitions is a consequence of the method that allows for the utilization of multiple reducers. One research direction, therefore, is to develop a scheme that balances the *reducer groups* in MR-GPMRS such that the computations is balanced across the reducers and the communication cost is minimized at the same time.

As the results show, using multiple reducers is not the best option when the skyline percentage is low. So for the algorithm to perform optimally for any data set, it is necessary to develop a scheme that allows MR-GPMRS to intelligently decide how many reducers to use. It is possible that

the method for choosing the PPD in the proposed algorithm can be improved.

The method does not incur a noticeable amount of extra runtime, and it chose a viable PPD for all the data sets tested for. It did not necessarily choose the optimal PPD, however, so finding a method that is guaranteed to choose the optimal PPD for any data set would be a significant improvement.

8. REFERENCES

- [1] Amazon. Amazon elastic compute cloud. <http://aws.amazon.com/ec2/>.
- [2] Apache. Welcome to apacheTMhadoop®! <http://http://hadoop.apache.org//>.
- [3] S. Brznsyi, D. Kossmann, and K. Stocker. The skyline operator. *17th International Conference on Data Engineering*, 2001.
- [4] L. Chen, K. Hwang, and J. Wu. Mapreduce skyline query processing with a new angular partitioning approach. *26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum*, 2012.
- [5] B. Cui, H. Lu, Q. Xu, L. Chen, Y. Dai, and Y. C. Zhou. Parallel distributed processing of constrained skyline queries by filtering. *24th ICDE*, 2008.
- [6] J. Dean and S. Ghemawat. Mapreduce: Simplified data processing on large clusters. *OSDI*, 2004.
- [7] Facebook. Hive - a petabyte scale data warehouse using hadoop. https://www.facebook.com/note.php?note_id=89508453919.
- [8] LinkedIn. A professional network built with java technologies and agile practices. <http://www.slideshare.net/linkedin/linkedins-communication-architecture>.
- [9] Twitter. Hadoop at twitter. <http://engineering.twitter.com/2010/04/hadoop-at-twitter.html>.
- [10] A. Vlachou, C. Doukeridis, and Y. Kotidis. Angle-based space partitioning for efficient parallel skyline computation. *SIGMOD '08*, 2008.
- [11] Yahoo. Hadoop at yahoo! <http://developer.yahoo.com/hadoop/>.
- [12] B. Zhang, S. Zhou, and J. Guan. Adapting skyline computation to the mapreduce framework: Algorithms and experiments. *DASFAA Workshops*, 2011.

Chapter 3

Implementation

3.1 MR-GPSRS and MR-GPMRS

The purpose of this section is to describe how the two algorithms Grid Partitioning Single Reducers Skyline (GPSRS) and Grid Partitioning Multiple Reducers Skyline (GPMRS) are implemented in Java for the Hadoop MapReduce framework.

3.1.1 Bitstring Generation

In this section, the implementation of the *bitstring generation* phase is described. The phase consists of a single MapReduce task consisting of multiple mappers and a single reducer. The mappers construct bitstrings BS_1 , BS_2 , ..., BS_i of subsets $R_1 \cup R_2 \cup \dots \cup R_n = R$ where R is the data set being processed. The mappers output their bitstrings to a single reducer that concatenates them into a single bitstring BS_R , i.e. $BS_1 \vee BS_2 \vee \dots \vee BS_i = BS_R$ representing empty and non-empty partitions of the R .

3.1.1.1 Mapper

The *map* function (Listing 3.1) of the bitstring generation phase translates the coordinates of a tuple into an index of its partition in the bitstring (line 3). The *Map* function does this for a number of different grid configurations in order for the *reducer* to make the best choice among several different grid configurations (lines 2-4)

Listing 3.1: Map function of the bitstring generation phase

```
1 public void map(LongWritable key, Text value, Context context)
2     ... {
3     for (int i = 0; i < cBitSet.length; i++) {
4         cBitSet[i].set(value.toString());
5     }
6     cLocalPoints++;
```

3.1.1.2 Reducer

The *Reduce* function of the bitstring generation phase takes a set of bitstrings BS_1, BS_2, \dots, BS_i from the mappers $map_1, map_2, \dots, map_i$. The set of the bitstrings from the mappers are then to be merged, e.g. a grid configuration $bs_i \in BS_1$ and $bs_i \in BS_2$ is to be merged into a single bitstring. The reducer then prunes the merged bitstrings $bs_{R.1}, bs_{R.2}, \dots, bs_{R.i}$ for dominated partitions. The pruning is done by iterating through the set bits (partitions) of the bitset and getting the indexes of dominated partitions (lines 5-11 of Listing 3.2). The partitions at the edge of the grid is not able to prune any other partitions and they do not have to be processed (line 6). When the iteration is done, every dominated partition is marked with a 1 in the bitstring *tAntiDominatingOffsets*, and by executing an AND NOT operation on the bitstring *cBitString* the dominated partitions are pruned (line 13). Finally, the reducer chooses the grid configuration that is the closest to the optimal average points per partition and outputs it.

Listing 3.2: Pruning

```
1 public void prune() {
2     BitSet tempBitSet = (BitSet) cBitString.clone();
3     BitSet tAntiDominatingOffsets = new BitSet();
4
5     for (int tOffset = tempBitSet.nextSetBit(0); tOffset >= 0;
6         tOffset = tempBitSet.nextSetBit(tOffset + 1)) {
7         if (isAtEdge(tOffset))
8             continue;
9         getDominatingOffsets(tOffset, tAntiDominatingOffsets);
10        tempBitSet.andNot(tAntiDominatingOffsets);
11    }
12
13    cBitString.andNot(tAntiDominatingOffsets);
14 }
```

3.1.2 Skyline Computation

In this section, the implementation of the *skyline computation* phase is described. The phase consists of a single MapReduce task consisting of multiple mappers and reducers. The mappers are initialized with the bitstrings which they use for pruning tuples in dominated partitions and pruning comparisons between tuples in partitions. The bitstring is also used to decide which reducers the partitions are sent to. The reducers merge the partitions from different mappers, and calculates its part of the skyline.

3.1.2.1 Mapper

The *Map* function of the skyline computation phase inserts each tuple into its corresponding partition. The tuple t is inserted into a partition p by calling the **BitList** method *insert(t)* (Listing 3.3). The index of the tuples partition is computed and used to check against the bitstring to see if it belongs to a pruned partition, discarding it if it does (lines 2-5). The tuple is then inserted into its partition which is stored in a *HashMap* with the index of the partition in the bitstring as a key (lines 7-10).

Listing 3.3: Map function of skyline computation

```
1 public void insert(String pTuple) {
2     int tOffset = getOffset(pTuple);
3
4     if (!cBitString.get(tOffset))
5         return;
6
7     if (!cPartitionList.containsKey(tOffset))
8         cPartitionList.put(tOffset, new Partition(tOffset));
9
10    cPartitionList.get(tOffset).insert(new Tuple(pTuple));
11 }
```

The mapper, however, cannot begin to output before every tuple has been sorted into its partition and the partitions have had their tuples pruned against the tuples of the partitions in their anti-dominating regions. The pruning of tuples in partitions is done by looping through every partition in the *HashMap* comparing them with any other partition in the *HashMap* which is in their anti-dominating region (lines 2-12 of Listing 3.4).

Listing 3.4: computeSkyline()

```
1 public void computeSkyline() {
2     for (Partition iPartition : cPartitionList.values()) {
3         BitSet antiDominating = getAntiDominatingOffsets(iPartition
4             .getcOffset());
5         antiDominating.clear(iPartition.getcOffset());
6
7         for (int i = antiDominating.nextSetBit(0); i >= 0; i =
8             antiDominating.nextSetBit(i + 1)) {
9             if (cPartitionList.containsKey(i)) {
10                Skyline.debugDominanceChecksBetweenPartitions++;
11                cPartitionList.get(i).check(iPartition);
12            }
13        }
14    }
```

The *independent groups* are found by iterating through a reverse copy of the global bitstring (lines 1-3 of Listing 3.5) and taking the first parti-

tion and its anti-dominating partitions as an *independent group* (line 4-5). The independent group is then cleared from the bitstring and the iterations continue until it is empty (line 6).

Listing 3.5: Independent Groups

```
1 BitSet pTempBitString = (BitSet) cBitSet.getBitString().clone()
  ;
2 PartitionHandler pPartitionHandler = new PartitionHandler(
  cNumberOfReducers);
3 for (int i = pTempBitString.length(); (i = pTempBitString.
  previousSetBit(i - 1)) >= 0;) {
4   BitSet cPartitions = cBitSet.getAntiDominatingOffsets(i);
5   pPartitionHandler.add(cPartitions);
6   pTempBitString.andNot(cPartitions);
7 }
```

The *reducer groups* are then constructed from the *independent groups* where the number of reducer groups is decided by the number of reducers as specified by a parameter value provided by the user. When the number of *reducer groups* is one, the algorithm is GPSRS. When the number of *reducer groups* is more than one, the algorithm is GPMRS.

However, in order for the mappers to output these *reducer groups*, responsibility of every partition has to be given to a *reducer group* (lines 2 and 5 of Listing 3.6). A *reducer group* requires an increased number of comparisons (partition comparisons) each time it has been made responsible for a partition. The *reducer groups* have to be balanced in terms of comparisons to balance the global skyline computation among the reducers. To balance the *reducer groups*, the group with the minimum number of computations is chosen to take responsible of a partition (lines 3-4). The *reducer groups* then take the non-responsible partitions that are within the anti-dominating region of the partitions they are responsible for (lines 8-10)

Listing 3.6: Assigning Responsibility

```
1 public void getPartitions(BitList cBitSet) {
2   while (cBitString.cardinality() != 0) {
3     ReducerGroup minReducerGroup = getMinimumGroup(cBitSet);
4     int responsiblePartition = minReducerGroup.takeResponsible(
      cBitSet);
5     cBitString.clear(responsiblePartition);
6   }
7
8   for (ReducerGroup item : cCollection) {
9     item.getNonResponsible(cBitSet);
10  }
11 }
```

The mappers serialize the *reducer group* objects and outputs them to the reducers.

3.1.2.2 Reducer

The purpose of the *Reduce* function of the skyline computation phase is to compute the global skyline of partitions its *reducer groups* are responsible for.

It begins when it has received the output from all its mappers. Since a reducer receives output from multiple mappers it is necessary to merge identical partitions, i.e. partitions describing the same space from different mappers (lines 9 and 13 of Listing 3.7). The merging is done such that only the tuples belonging to the partitions local skyline is maintained. The responsible and non-responsible partitions are kept separate and *computeSkyline()* is called on the responsible partitions while the non-responsible is used to prune the responsible partition to ensure that the tuples in the responsible partitions are a part of the global skyline.

Listing 3.7: Merging Inputs

```
1 public void reduce(IntWritable key, Iterable<BytesWritable>
   values, Context context) ... {
2     ...
3     for (BytesWritable item : values) {
4         ...
5         ReducerGroup reduceGroup;
6         reduceGroup = kryo.readObject(input, ReducerGroup.class);
7         ...
8         for (Partition partition : reduceGroup.responsible()) {
9             cResponsible.merge(partition);
10        }
11
12        for (Partition partition : reduceGroup.nonResponsible()) {
13            cNonResponsible.merge(partition);
14        }
15    }
16    ...
```

The reducer then outputs the tuples of the responsible partitions as the result.

3.2 MR-BNL

The purpose of this section is to describe how the algorithm MapReduce - Block Nested Loop (MR-BNL) from [12] is implemented in Java for the Hadoop MapReduce framework.

3.2.1 Phase 1

In this section, the implementation of the first phase of the MR-BNL algorithm is described. The purpose of the first phase is to calculate the local

skyline of the subspaces generated by utilizing a grid partition scheme with two subspaces per dimension, i.e. 2^d subspaces where d is the number of dimensions.

3.2.1.1 Mapper

The purpose of the mappers in the first phase is to partition the dataset to multiple reducers based on their subspace flag, e.g. 01, 00, 10, and 11 in case of a 2 dimensional data set. The subspaces are split into *higher* (1) and *lower* (0) based on the mean of the data range. In Listing 3.8 the tuples are read and partitioned according to their subspace (lines 3-5)

Listing 3.8: Phase 1 - Map

```
1 public void map(LongWritable key, Text value, Context context)
2     ... {
3     Tuple tuple = new Tuple(value);
4     keyOut.set(tuple.getFlagString());
5     context.write(keyOut, value);
6 }
```

3.2.1.2 Reducer

The *Reduce* function (Listing 3.9) calculates its local skyline of its assigned subspace, which have been partitioned by the mapper. This is done by keeping a window of skyline points (line 3-5) and inserting new skyline tuples while removing dominated tuples from the window. The window itself is a self organising list which moves frequent dominating tuples to the top. The reason for this is that less iterations through the list has to be made if the most dominating tuples are at the top.

Listing 3.9: Phase 1 - Reduce

```
1 public void reduce(Text key, Iterable<Text> values, Context
2     context) ... {
3     Window cSkylineWindow = new Window();
4     for (Text value : values) {
5         cSkylineWindow.judge(new Tuple(value));
6     }
7     for (Tuple tuple : cSkylineWindow) {
8         context.write(new IntWritable(1), new Text(tuple.toString()
9             ));
10    }
```

3.2.2 Phase 2

In the following section, the implementation of the second phase of MR-BNL algorithm is described. The purpose of this phase is to calculate the global skyline from the set of local skylines produced by the reducers in the previous phase.

3.2.2.1 Mapper

The global skyline has to be calculated in a single task (map task or reduce task), as the algorithm needs to be aware of all the points to compute the global skyline. This algorithm computes its global skyline in a single reducer, which is done by having the mappers output the local skylines of the previous phase to a single reducer.

3.2.2.2 Reducer

The reducer partitions the tuples based on their subspace flag, while keeping the local skyline of the partition updated. The tuples belonging to the dominated partition, i.e. the partition *higher* in every dimension, is discarded if there is tuples in the dominating partition, i.e. the partition *lower* in every dimension. In Listing 3.10, the partitioning of the tuples is done by keeping a window for each partition in a **HashMap** (line 4-8).

Listing 3.10: Phase 2 - Reduce

```
1 public void reduce(IntWritable key, Iterable<Text> values,
2   Context context) ... {
3   for (Text item : values) {
4     Tuple newTuple = new Tuple(item);
5     if (!cPartitions.containsKey(newTuple.getFlagString())) {
6       cPartitions.put(newTuple.getFlagString(),
7         new Window(newTuple.getFlag()));
8     }
9     cPartitions.get(newTuple.getFlagString()).judge(newTuple);
10  }
11  ...
```

The global skyline is computed by comparing comparable partitions and pruning dominated points from each partition (lines 2-5 of Listing 3.11). The reducer is then able to output the final skyline.

Listing 3.11: Phase 2 - Skyline Computation

```
1 public void insert(Window value) {
2   for (Entry<String, Window> mWindow : cWindow.entrySet()) {
3     if (isComparable(mWindow.getValue().getFlag(), value.
4       getFlag()))
5       mWindow.getValue().check(value);
```



```
5    }  
6  
7    cWindow.put (Arrays.toString(value.getFlag()), value);  
8 }
```

3.3 MR-ANGLE

The purpose of this section is to describe how the algorithm MapReduce - Angular (MR-Angular) from [4] is implemented in Java for the Hadoop MapReduce framework.

3.3.1 Phase 1

In this section, the implementation of the first phase of the MR-Angular algorithm is described. The purpose of the first phase is to partition the data set according to their angular coordinates and compute the local skyline of the angular partition.

3.3.1.1 Map

The *Map* function translates the cartesian coordinates of the tuples into n-sphere coordinates, using the formula in [10], and partitions the tuples based on the equi-volume partitioning scheme as described in [10] (lines 2-4 of Listing 3.12).

Listing 3.12: Map of first phase

```
1 public void map(LongWritable key, Text value, Context context)  
2     ... {  
3     Tuple vT = new Tuple(value.toString());  
4     keyOut.set(vT.getPartition());  
5     context.write(keyOut, value);  
6 }
```

The implementation of the formula is shown in Listing 3.13 where a coordinate $x = [x_1, x_2, ..x_d]$ is iterated through using x_1 to x_{d-1} as the divisor (lines 3). The dividend tC is constructed by iterating backwards through the cartesian coordinates starting with x_d and ending with x_i , the current dividend (lines 5-15). A special case is when calculating the last n-sphere coordinate where the translation is done differently (lines 7-11 and 17-21).

Listing 3.13: Translation of coordinates

```
1 public int getPartition() {  
2     Double[] tNSphereCoord = new Double[cTuple.length - 1];  
3     for (int i = 0; i < cTuple.length - 1; i++) {
```

```
4     double tC = 0;
5     for (int j = cTuple.length - 1; j > i; j--) {
6         if (j == cTuple.length - 1) {
7             if (i == cTuple.length - 2) {
8                 tC += cTuple[j];
9             } else {
10                tC += Math.pow(cTuple[j], 2);
11            }
12        } else {
13            tC += Math.pow(cTuple[j], 2);
14        }
15    }
16
17    if (i == cTuple.length - 2) {
18        tC = tC / cTuple[i];
19    } else {
20        tC = (Math.sqrt(tC)) / cTuple[i];
21    }
22    tNSphereCoord[i] = Math.atan(tC);
23 }
24 ...
```

3.3.1.2 Reduce

The *Reduce* function is responsible for calculating the skyline of a single partition. The skyline computation is implemented by using BNL with a window and a self organizing list as described in [4]. The resulting skyline tuples are then output.

3.3.2 Phase 2

In the following section, the implementation of the second phase of the MR-Angular algorithm is described. The purpose of the second phase is to compute the global skyline by having the mappers read the local skylines from the previous phase and output it to a single reducer.

3.3.2.1 Map

As for MR-BNL, this algorithm computes its global skyline in a single reducer by having the mappers output the local skylines of the previous phase to a single reducer.

3.3.2.2 Reduce

The global skyline computation is performed like the local skyline computation in the reducers of the first phase. The *Reduce* function takes the tuples from the mappers as input and calculates the global skyline using BNL.

Appendix A

CD-ROM

The CD-ROM contains the following:

- A folder called Report containing the project report.
- A folder called Algorithms containing the source code for the algorithms.