

**SCALABLE MACHINE LEARNING ALGORITHMS IN  
PARALLEL DATABASE SYSTEMS EXPLOITING A  
DATA SUMMARIZATION MATRIX**

---

A Thesis

Presented to

the Faculty of the Department of Computer Science

University of Houston

---

In Partial Fulfillment

of the Requirements for the Degree

Master of Science

---

By

Yiqun Zhang

May 2016

**SCALABLE MACHINE LEARNING ALGORITHMS IN  
PARALLEL DATABASE SYSTEMS EXPLOITING A  
DATA SUMMARIZATION MATRIX**

---

Yiqun Zhang

APPROVED:

---

Carlos Ordonez, Ph.D., Chairman  
Dept. of Computer Science

---

Omprakash Gnawali, Ph.D.  
Dept. of Computer Science

---

Xin Fu, Ph.D.  
Dept. of Electrical & Computer Engineering

---

Dean, College of Natural Sciences and Mathematics

# Acknowledgements

I would like to offer my sincere thanks to my advisor Dr. Carlos Ordonez for his patience, guidance and support during my study at the University of Houston. His knowledge, insights and generous help gave me great confidence in my efforts to undertake the challenge of the research project this thesis presents.

I am very grateful to Dr. Omprakash Gnawali and Dr. Xin Fu for taking the time to serve on my thesis committee and sharing their knowledge and experience. It is my privilege to work with my fellow researchers, Wellington Cabrera and David Matusevich, I want to thank them for the friendship and teamwork we had together during the past few years.

I also want to thank my parents Jingli and Desheng, and other family members especially my cousin Yuanfeng Wen and my sister-in-law Wensi Song. Their unconditional love has always been the strongest emotional support that I am so blessed to have through my graduate study here, more than 7,000 miles away from home, in the United States.

My sincere thanks also go to my dear friend Xi Liu for his company and support, as well as Can Cao, Chen Shen, Yupeng Zhang, Ziyi Liu, Tao Feng, Qiang Li, Fangkai Yang, Zihao Zhou, Yin Xu, Yifei Wan, and Qinqin Lyu.

**SCALABLE MACHINE LEARNING ALGORITHMS IN  
PARALLEL DATABASE SYSTEMS EXPLOITING A  
DATA SUMMARIZATION MATRIX**

---

An Abstract of a Thesis

Presented to

the Faculty of the Department of Computer Science

University of Houston

---

In Partial Fulfillment

of the Requirements for the Degree

Master of Science

---

By

Yiqun Zhang

May 2016

# Abstract

Data summarization is an essential mechanism to accelerate analytic algorithms on large data sets. In this work we present a comprehensive data summarization matrix, namely the Gamma matrix, from which we can derive equivalent equations for many analytic algorithms. In this way, iterative algorithms are changed to work in two phases: (1) Incremental and parallel summarization of the data set in one pass; (2) Iteration in main memory exploiting the summarization matrix in many intermediate computations. We show our summarization matrix captures essential statistical properties of the data set and it allows iterative algorithms to work a lot faster in main memory. Specifically, we show our summarization matrix benefits statistical models, including PCA, linear regression and variable selection. From a system perspective, we carefully study the efficient computation of the summarization matrix in two parallel database systems including the array DBMS SciDB, and the columnar relational DBMS HP Vertica. We also propose general optimizations according to the data density and system-dependent optimizations for each platform. We present an experimental evaluation benchmarking system and algorithm performance. Our experiments show that our algorithms work significantly faster than existing machine learning libraries for model computations in R and Spark, and R working together with SciDB in general can run our algorithm significantly faster than all the other parallel analytic systems compared. More importantly, it eliminates main memory and performance limitations from R.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Work</b>	<b>5</b>
2.1	Faster Analytic Algorithms . . . . .	5
2.2	Data Summarization . . . . .	6
2.3	Scalable Matrix Multiplication . . . . .	8
<b>3</b>	<b>Definitions</b>	<b>10</b>
<b>4</b>	<b>Defining and Exploiting a Multidimensional Summarization Matrix</b>	<b>11</b>
4.1	Gamma Summarization Matrix . . . . .	11
4.2	Two-Phase Analytic Algorithm . . . . .	16
4.3	Models Exploiting the Gamma Matrix . . . . .	17
4.3.1	Principal Component Analysis (PCA) . . . . .	18
4.3.2	Linear Regression (LR) . . . . .	20
4.3.3	Variable Selection (VS) . . . . .	21
<b>5</b>	<b>Data Summarization in Parallel Database Systems</b>	<b>25</b>
5.1	Data summarization in SciDB . . . . .	25
5.1.1	Programming Mechanisms in SciDB . . . . .	26
5.1.2	The Gamma Array Operator to Summarize a Matrix . . . . .	26
5.1.3	Additional Considerations . . . . .	32
5.2	Data Summarization in Columnar DBMS . . . . .	35
5.3	Time Complexity and Parallel Speedup . . . . .	38
<b>6</b>	<b>Experimental Evaluation</b>	<b>41</b>

6.1	Experimental Setup . . . . .	42
6.1.1	Mathematical Models Functions and Parameters . . . . .	43
6.1.2	System Setup and Tuning . . . . .	43
6.2	Comparing Analytic Systems ( $N = 1$ node) . . . . .	44
6.2.1	Comparing R+SciDB and R alone . . . . .	44
6.2.2	Computing $\Gamma$ on R, Array DBMS and Column DBMS . . . . .	47
6.2.3	Comparing $\Gamma$ Algorithm with Fast Matrix Multiplication: SciDB versus Intel MKL (LAPACK) . . . . .	48
6.3	Parallel Processing ( $N = 1, 10, 100$ nodes) . . . . .	51
<b>7</b>	<b>Conclusions</b>	<b>55</b>
	<b>Bibliography</b>	<b>57</b>

# List of Figures

4.1	The workflow for a system using $\Gamma$ for model computation . . . . .	16
5.1	How the Gamma operator works in SciDB . . . . .	34
5.2	The horizontal and vertical layout for $X$ in the columnar DBMS. . .	36
6.1	Comparing systems to compute $\Gamma$ on local server: Array DBMS, R, SQL. . . . .	49



# List of Tables

4.1	Summary of matrices: summarization and models. . . . .	15
6.1	Data sets. . . . .	42
6.2	Hardware and Operating System. . . . .	42
6.3	Comparing computation of model $\Theta$ using R and DBMS+R; dense matrix operator; data set KDDnet; local server; times in secs. . . . .	46
6.4	Bayesian variable selection (VS) in gene data set ( $d = 12506, n = 248$ ): Comparing R+SciDB with R at 1000 iterations; times in secs. . . . .	46
6.5	Comparing matrix summarization algorithms; data set KDDnet; dense ( $dn$ entries) and sparse (zeroes deleted); Local 1/2; times in secs. . . . .	48
6.6	Comparing dense operator (working on disk+RAM), sparse operator (working on disk+RAM) and dense matrix multiplication in MKL BLAS (parallel LAPACK, working only on RAM) to compute $\Gamma$ ; local server 4GB RAM; synthetic data sets; times in secs. . . . .	50
6.7	Parallel speedup: $\Gamma$ dense matrix summarization operator with multi-threaded processing varying $M = \#$ SciDB instances ( $N = 1$ node, 4-core CPU); times in secs. . . . .	52
6.8	Parallel processing: Comparing R+SciDB and Spark varing $N$ ( $\#$ of nodes); data set KDDnet $d = 38$ ; dense matrix; times in secs. . . . .	52

# Chapter 1

## Introduction

Business and scientific activities are becoming more and more data-driven. Most organizations are now collecting data at the most detailed level possible and the volume of data is growing very fast at ever-accelerating rates. As a result, the increasing need for analyzing and understanding collected data is drawing more and more people's attention to the great challenges of big data analytics.

Big data analytics is a new data mining trend [13, 1] with higher data volume, varied data types, and parallel processing. From a system prospective, row DBMSs [30] remain the best technology for transaction processing [29] and column DBMSs are becoming a competitor in query processing [26] in large databases [27]. However, both kinds of DBMS are considered slow and difficult to use for mathematical computations on large data sets due to the abundance of matrices. Currently, Hadoop and MapReduce systems [15, 27] (e.g. Spark [34, 31], HadoopDB, Cassandra, Hive, MongoDB, the Mahout project, and so on) are the most popular for big data analytics. From a statistical perspective R is the most popular platform, but despite many efforts it remains slow and not scalable for big data analytics. Array DBMSs (SciDB [28], Rasdaman [4]) represent alternative analytic systems, closer to relational DBMSs, that enable computation with large matrices in R, which share many features with old DBMSs like efficient I/O, indexing, concurrency control and parallel processing, but have significantly different storage and query processing mechanisms.

On the algorithm side, there has been extensive research efforts on optimizing machine learning algorithms on large-scale data sets and in parallel. Many algorithms are adapted to work iteratively or using a randomized approach. However, those work still require scanning the original data set multiple times. As the data sets grow explosively, those algorithms can suffer degraded performance due to the I/O bottleneck. In contrast, optimizing machine learning algorithms with data summarization techniques has received scant attention. In this work, we focus on the optimization of algorithms to compute three fundamental and complementary models. Specifically, we study the computation of principal component analysis (PCA) [14] with the correlation matrix [32], linear regression (LR) [14] and variable selection (VS) on large data sets. However, these three models require different numerical and statistical methods for their solutions: PCA is generally computed with a Singular Value Decomposition (SVD) on the correlation [32, 14] (covariance) matrix of the data set; LR is solved with least squares via a matrix factorization; VS requires visiting an exponential search space of variable combinations, whose most promising solution is currently found by MCMC methods, used in Bayesian statistics [11]. A costly matrix multiplication plays a very important role in the computation of those models. The root causes are a matrix transposition and then a slow matching between the  $i$ th row from the input matrix with the  $j$ th column from its transpose. With that motivation in mind, we introduce a summarization matrix, called Gamma, that captures essential statistical properties of the data set, to compute the first and the second moment (mean, variance) of multivariate probability distributions. The most important characteristics of the Gamma matrix is that it is much smaller than the data set and it can be used to avoid recomputing costly matrix multiplications. We show the Gamma summarization matrix can be efficiently computed via a novel

form of matrix multiplication, optimized for large data sets, in which Gamma can be incrementally updated in main memory. To that end, we introduce two parallel incremental algorithms for dense and sparse large input matrices (i.e., the data set), respectively, which scan the data set once, in parallel. Exploiting Gamma, we introduce a “meta” algorithm that divides model computation of our models into two phases: (1) parallel incremental summarization to get Gamma, (2) model computation in main memory exploiting the Gamma matrix, including the iterative behavior.

Experimental evaluation shows our summarization operator removes main memory limitations from R and it is much faster than R (when the data set fits in RAM) and orders of magnitude faster than a column DBMS (one of the fastest DBMSs, evaluating summarization with SQL queries). There are efficient functions in R, LAPACK [3, 8], SciDB [28] and Spark [34] to compute matrix multiplication, but we show they are inefficient to summarize the data set. More importantly, we show R and LAPACK fail when there is insufficient RAM for large input matrices. From a parallel perspective, the array DBMS is orders of magnitude faster and has better speedup than Spark, when running in a large cluster in the cloud.

We discuss closely related work in Chapter 2. Chapter 3 is a reference chapter that introduces definitions used throughout the thesis. Chapter 4 presents our theory and research contributions, divided into two main aspects: (1) a comprehensive summarization matrix that benefits several linear models that allow statistical methods to work in two phases: summarization and iteration. (2) Deriving equivalent matrix equations based on summarization that decrease iterations or the number of times the data set must be read. Chapter 5 proposes the system implementations for the

summarization computation that can work in parallel in one pass on the data set, and the optimization details in every system. We also explain how the summarization matrix is exploited in the R package, taking advantage of its powerful mathematical capabilities. Chapter 6 presents an experimental evaluation comparing our solution to state-of-the-art analytic systems including the R package, Vertica, and Spark. We also specifically compare matrix multiplication with LAPACK, the standard in HPC. Conclusions and directions for future work are discussed in Chapter 7.

# Chapter 2

## Related Work

This thesis is established on several our previous research achievements. Our summarization matrix was first proposed in [24] to solve PCA, linear regression, and variable selection. Our recent research [25] explored the optimizations for sparse input data sets and showed that our SciDB implementation defeated Spark in performance on a large cluster on the cloud. In [36] we migrated our algorithm to a relational columnar DBMS and its performance is very close to that of SciDB.

### 2.1 Faster Analytic Algorithms

Research has developed fast algorithms based mostly on sampling, data summarization, and gradient descent [12, 16], but generally working in a sequential manner (data mining) or outside a parallel data analysis system (Hadoop or DBMS). Sample-based algorithms [9] require mechanisms to control approximation error and they are hard to program in parallel for data sets with skewed distributions. In our case, no sampling is required, although it can further accelerate our summarization algorithms. Stochastic (incremental) gradient descent (SGD) [15] is another popular approach, useful when there is a convex function to optimize (like least squares in LR). As drawbacks, SGD is naturally sequential (difficult to process in parallel), it obtains an approximate solution and it is difficult to adapt to non-convex functions (e.g. clustering). On the other hand, we believe we are the first to study parallel

data summarization as matrix multiplication. Given their relevance to our work, data summarization and matrix multiplication are discussed in more detail below.

Despite the prominence of statistical and machine learning models, there is not much work studying how to accelerate their computation inside a DBMS. Accelerating SVD has received attention [20], but producing approximate solutions and without considering parallelism. Linear classification with Support Vector Machines (SVMs) is not a problem we considered, but it is straightforward to define a summarization operator for a diagonal matrix  $\Gamma$  to get a reasonably accurate Naive Bayes classifier. The importance of aggregate UDFs to accelerate the computation of statistical models in relational DBMSs is identified in [23], making a big step forward compared to a pure query-based approach [21]. Following the same ideas, [15] introduces the MADlib library to compute statistical models with SQL mechanisms, combining queries and UDFs. Nevertheless, neither [23] nor [15] had envisioned matrix multiplication as a vector-based outer product computation which can be pushed into an aggregation operator with arrays both as input and output. Moreover, the integration with R was not considered.

## 2.2 Data Summarization

A similar, but less general, data summarization to ours was pioneered in [35] to accelerate the computation of distance-based clustering: the sums of values and the sums of squares. Later [5] exploited such summaries as multidimensional sufficient statistics for the K-means and EM clustering algorithms. The main differences with

[35] and [5] are: data summaries were useful only for one model (clustering). Compared to our proposed matrix, their summaries represent a (constrained) diagonal version of  $\Gamma$  because dimension independence is assumed (i.e., cross-products, covariances, correlations, are ignored) and there is a separate vector to capture  $L$ . From a computational perspective, our summarization algorithm boils down to one matrix multiplication, whereas those algorithms work are aggregations. Another major difference is that in our models one summarization matrix is sufficient, whereas those clustering models need more than one matrix. In short, our summarization is more general and it can help computing more complex models like PCA, LR and VS that could not be solved with older summaries. From a statistical perspective, we have identified the product of the extended matrix  $Z$  with itself as a fundamental summarization for a data set covering zeroth, first and second moments of its probabilistic distribution. A more general data summarization capturing up to the fourth moment was proposed in [10], but it relies on binning (i.e., building histograms) which are incompatible with most statistical methods. Parallel processing for data summarization has received moderate attention. Reference [18] highlights the following techniques: sampling, incremental aggregation, matrix factorization and similarity joins. Our proposal is a combination of incremental aggregation and scalable matrix multiplication that enables fast matrix factorization in main memory. A closely related work that identified the linear sum of points  $L$  and the quadratic sum of points  $Q$  with cross-products is [23], but it did not recognize the importance of expressing the computation as a single matrix product, did not study parallel processing, did not consider sparse matrices and did not synthesize the 2-step algorithm, which has potential for new models. From a “systems” angle, array DBMSs supporting UDFs did not exist and it was not envisioned the summarization step could be solved entirely



by the DBMS and the iterative and mathematical computations being evaluated in R.

## 2.3 Scalable Matrix Multiplication

Matrix multiplication has been extensively studied with dense and sparse matrices, as well as one processor (sequentially) or  $N$  processors (in parallel) [3, 7]: parallel sparse matrix multiplication is the hardest combination. The specific assumptions about matrix shape, density and parallel computation model vary widely. Most research has proposed algorithms partitioning and storing the input matrices by block in distributed memory [7, 8]. Dense matrix multiplication in main memory in one computer is considered a solved problem for the most part [7]. Even though multiplication of large matrices exceeding RAM limits (out-of core) has been studied before [7], nobody has looked at the case that the input matrix (data set) is large, but the result matrix is much smaller and dense. In other words, most research has focused on dense-dense multiplication producing a dense matrix or sparse-sparse multiplication producing a sparse matrix. From a parallel perspective, most algorithms work in a multi-core CPU or GPU. In consequence, parallel matrix multiplication algorithms for matrices residing on secondary storage have not made their way into ScaLAPACK because of the complexity of combining parallel computation with MPI and efficient I/O on disk. MPI is not an efficient interface for DBMS technology because it has an underlying shared-memory model that requires transmitting data across nodes.

Sparse and parallel matrix multiplication are harder problems [2, 6, 17, 33], where the parallel computation model is shared memory via MPI. That is, such approaches

are generally incompatible with shared-nothing architectures like Hadoop and parallel DBMSs. However, there is some work from the database systems angle to multiply sparse matrices for graph analytics [22]. As additional differences, most algorithms assume matrices of arbitrary shape and size, as long as they are compatible for matrix multiplication and they reduce the problem to a matrix-vector multiplication. In contrast, our matrix summarization operator is significantly different (but useful in big data analytics), because it is optimized for a single rectangular matrix as input (in general the shape of a large data set), the input matrix is block-partitioned only along the largest dimension ( $n$ ), parallelization is based on vector-vector outer products and there is no communication overhead during parallel processing. Last but not least, an array DBMS provides not only fast mechanisms to manipulate large matrices, but also an elegant and well-defined programming mechanism to create new matrix operators, something that was not feasible with SQL-based DBMSs.

# Chapter 3

## Definitions

We start by defining the input matrix, stressing it is a set of  $n$  *column* vectors. All models take a  $d \times n$  matrix  $X$  as input. Let  $X = \{x_1, \dots, x_n\}$  be the input data set with  $n$  points, where each point  $x_i$  is a vector in  $\mathbf{R}^d$ . Intuitively,  $X$  is a wide rectangular matrix. In the case of LR and VS,  $X$  is augmented with a  $(d + 1)$ th dimension containing an output variable  $Y$ , making  $X$  a  $(d + 1) \times n$  matrix. We use  $i = 1 \dots n$  and  $j = 1 \dots d$  as matrix subscripts. We emphasize that for convenience in mathematical notation we use *column* vectors and *column*-oriented matrices.

We use  $\Theta$ , a set of matrices and associated statistics, to refer to a statistical model in a generic manner. Thus  $\Theta$  represents the PCA, LR, and VS models defined in Chapter 4.3. Moreover,  $\Theta$  could represent a clustering or classification model. PCA represents an unsupervised (descriptive) model to reduce dimensionality. Linear regression is a fundamental supervised model, whose solution helps understanding and building other linear models. Variable selection (VS), a supervised model, is among the computationally most difficult problems in statistics and machine learning, in which Markov Chain Monte Carlo (MCMC) methods [14, 11] are a promising solution.

# Chapter 4

## Defining and Exploiting a

## Multidimensional Summarization Matrix

We start by introducing basic vectors and matrices containing sufficient statistics [14], which are generalized and integrated into a single summarization matrix  $\Gamma$  that is exploited in intermediate demanding computations for the models introduced in Chapter 3. Chapter 4.1 introduces  $\Gamma$  and explains its mathematical properties, which enable scalable parallel algorithms. Chapter 4.2 introduces a general algorithm to efficiently compute models in two phases. Chapter 4.3 explains in technical detail how  $\Gamma$  is exploited in intermediate computations. All “systems” aspects and parallel algorithms are presented in Chapter 5.

### 4.1 Gamma Summarization Matrix

In this section we introduce the  $\Gamma$  summarization matrix. We first explain this matrix contains several vectors and submatrices that represent important sums derived from the data set. In Chapter 4.3, we show such sums help deriving fundamental statistical properties of the data set, for each dimension (variable) and for each variable pair.

We first review sufficient statistics matrices [14, 23], which are integrated and generalized into a single matrix:

$$n = |X| \tag{4.1}$$

$$L = \sum_{i=1}^n x_i \tag{4.2}$$

$$Q = XX^T = \sum_{i=1}^n x_i \cdot x_i^T \tag{4.3}$$

Intuitively,  $n$  counts points,  $L$  is a linear sum of  $x_i$  and  $Q$  is a quadratic sum of  $x_i$ , where  $x_i$  multiplied by itself (i.e., “squared”) with a vector outer product. As explained in Section 4.3.2, the linear regression model uses an *augmented* matrix, represented by  $\mathbf{X}$ . We introduce a more general augmented matrix  $\mathbf{Z}$ , by appending an additional  $d + 1$ th row to  $\mathbf{X}$ , which contains the vector  $Y$ . Since  $X$  is  $d \times n$ ,  $\mathbf{Z}$  has  $(d + 2)$  rows and  $n$  columns, where row  $[0]$  are 1s and row  $[d + 1]$  is  $Y$ .

Matrix  $\Gamma$  contains a comprehensive, accurate and sufficient summary of  $X$  to efficiently compute all models previously defined. Below we show  $\Gamma$  in two equivalent forms: (1) as vector-matrix and matrix-matrix multiplications; (2) as sums of vector outer products. Notice  $\mathbf{1}$  is a column-vector of  $n$  1s, which allows expressing a sum as a matrix product. Such equivalence has important performance implications depending on how the matrix is processed.

$$\begin{aligned}\Gamma &= \begin{bmatrix} n & L^T & \mathbf{1}^T \cdot Y^T \\ L & Q & XY^T \\ Y \cdot \mathbf{1} & YX^T & YY^T \end{bmatrix} \\ &= \begin{bmatrix} n & \sum x_i^T & \sum y_i \\ \sum x_i & \sum x_i x_i^T & \sum x_i y_i \\ \sum y_i & \sum y_i x_i^T & \sum y_i^2 \end{bmatrix}\end{aligned}$$

The fundamental property of  $\Gamma$  is that it can be computed by a single matrix multiplication using  $\mathbf{Z}$ . Therefore, we study how to compute the matrix product below, related to, but not the same as the Gram matrix  $\mathbf{Z}^T \mathbf{Z}$  [7]:

$$\Gamma = \mathbf{Z}\mathbf{Z}^T \tag{4.4}$$

Matrix  $\Gamma$  is fundamentally the result of “squaring” matrix  $\mathbf{Z}$ .  $\Gamma$  is comparatively much smaller than  $X$  for large  $n$ , symmetric and computable via vector outer products. Such facts are summarized in the following properties:

*Property 1 (matrix product versus vector outer product):*  $\Gamma$  can be equivalently computed as follows:

$$\Gamma = \mathbf{Z}\mathbf{Z}^T = \sum_{i=1}^n z_i \cdot z_i^T. \tag{4.5}$$

*Property 2 (small size):* Given a large data set  $X$  where  $d \ll n$ , matrix  $\Gamma$  size is  $O(d^2) \ll O(dn)$  as  $n \rightarrow \infty$ .

*Property 3 (symmetry):*  $\Gamma = \Gamma^T$ , which can be condensed into a triangular matrix with  $d + (d + 2)^2/2 \approx d^2/2$  entries, saving 1/2 of CPU operations.

*Property 4 (first and second moment):* Matrix  $\Gamma$  summarizes  $X$  to compute the first and second moment of several multivariate probabilistic distributions.

Property 1 is fundamental to derive a highly efficient summarization operator. Property 1 gives two equivalent expressions to obtain  $\Gamma$ : one as a matrix multiplication and a second one as a sum of vector *outer products*. We make the following fundamental observation. Even though it is tempting to evaluate  $\mathbf{Z}\mathbf{Z}^T$  as a matrix multiplication it is a bad idea: we need to compute and materialize  $\mathbf{Z}^T$  first, a costly transposition, storing it and then perform the actual matrix multiplication. Instead, we defend the idea of evaluating the summation  $\sum_{i=1}^n z_i \cdot z_i^T$ . We will experimentally prove this hypothesis. Assuming  $z_i$  fits in main memory then it is reasonable to assume the vector outer product also fits in main memory. Therefore, it makes more sense to evaluate such sum in parallel. Given the prominent importance of  $\Gamma$  we believe our operator should be available on any big data analytic system. We call our operator  $\Gamma(X)$  (called `Gamma(X)` in a C++ function).

Property 2 means that if  $\Gamma$  can fit in main memory it is feasible to maintain a summary of  $X$  in main memory. However, we emphasize we cannot assume  $X$  can fit in main memory. Therefore, the challenge is to efficiently compute  $\Gamma$ , minimizing the number of times  $X$  must be read from secondary storage.

Property 3 states it is possible to save one half of CPU work. Notice storage space in triangular form is still  $O(d^2)$ , but it requires a nested indexing mechanism. In order to allocate a square array, which enables faster address computation and block-based access we prefer not to use compact storage. Notice Property 2 implies  $\mathbf{Z}\mathbf{Z}^T = (\mathbf{Z}\mathbf{Z}^T)^T = (\mathbf{Z}^T)^T \mathbf{Z}^T = \Gamma$ , highlighting the matrix product is non-commutative. In other words,  $\Gamma \neq \mathbf{Z}^T \mathbf{Z}$ . Table 4.1 provides a summary of all matrices, including the

Table 4.1: Summary of matrices: summarization and models.

Matrix	Size	Description
$X$	$d \times n$	Data set: dimensions/variables
$\mathbf{X}$	$(d + 1) \times n$	Augmented matrix with 1s
$\mathbf{Z}$	$(d + 2) \times n$	Augmented $X$ matrix with 1s and $Y$
$\Gamma$	$(d + 2) \times (d + 2)$	Summarization matrix
$U$	$d \times d$	Principal components
$V$	$d \times d$	Squared eigenvalues; diagonal
$Y$	$1 \times n$	Dependent variable
$\beta$	$(d + 1) \times 1$	Regression coefficients, $Y$ intercept
$\gamma$	$d \times 1$	Selected variables; binary vector

input data set, the output variable, model  $\Theta$  and  $\Gamma$ .

We explain Property 4 in more technical detail.  $\Gamma$  is a fundamental summarization matrix because it can help computing the first and second expected moments of many probabilistic distributions: the first moment  $E[x] = \mu$ , to get the global mean (average) and the second moment  $E[(x - \mu)(x - \mu)^T] = V$  for the covariance matrix, which measures mean squared error (MSE) per dimension and helps understanding each variable pair behavior. Mean squared error is the most commonly used error measure in statistics and machine learning due to being mathematically easier to optimize [14].  $\Gamma$  can be used to *directly* derive the global mean and the covariance matrix: That is,  $\Gamma$  helps understanding each individual variable and each pair of variables. The following models represent a generalization of these basic statistical matrices. Higher order variable interactions can be analyzed adding dimensions to  $X$  for each variable pair.



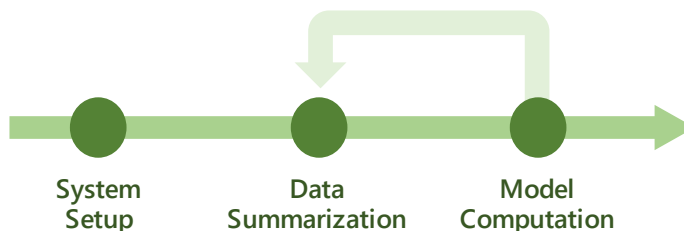


Figure 4.1: The workflow for a system using  $\Gamma$  for model computation

## 4.2 Two-Phase Analytic Algorithm

Recall from Chapter 3  $\Theta$  represents a statistical model. Therefore, based on  $\Gamma$  a fast algorithm to compute  $\Theta$  in two phases is the following:

1. Compute  $\Gamma$ .
2. Iterate exploiting  $\Gamma$  in intermediate matrix computations.

These two phases can be repeated as many times as needed for iterative algorithms in a system using this paradigm, workflow shown in Figure 4.1. This paper focuses on optimizing Phase 1, which is a well-defined problem for any input data set  $X$ . Phase 2 requires incorporating  $\Gamma$  in the steps of numerical and statistical methods. By exploiting  $\Gamma$  it becomes possible to reduce the number of times  $X$  is read, and to reduce CPU computations in iterative methods. Identifying in which models intermediate matrix computations accept  $\Gamma$  is a deep research topic. We point out that this works builds on top of previous research and we make clear for which kind of models our optimizations apply. That is, our choice of methods to solve PCA, LR and VS has been carefully done based on the fact that they can exploit  $\Gamma$  in the

most demanding matrix equations.

### 4.3 Models Exploiting the Gamma Matrix

Our summarization matrix benefits a big family of statistical models. In this thesis we focus on PCA, LR, and VS, where one summarization matrix is sufficient despite the fact that their iterative methods are different. Our summarization matrix can be generalized to compute the Naive Bayes classifier, K-means/EM clustering, logistic regression and Linear Discriminant Analysis, among others. But these latter models require more than one summarization matrix (e.g.,  $k$  summarization matrices for K-means or EM), which introduces different assumptions and other iterative methods. Therefore, such models are a research issue for future work.

For the statistical models we focus on we use a single summarization matrix  $\Gamma$  that is non-diagonal. That is, we do not assume dimensions are independent. Therefore,  $\Gamma$  is a full matrix with  $O(d^2)$  non-zero entries. From an algorithmic perspective, we have managed to use  $\Gamma$  as a proxy of  $X$  in every step. Therefore, the model  $\Theta$  is computed in one (parallel) pass over  $X$ , as explained in the next subsections.

There are two basic directions to generalize  $\Gamma$ : (1) computing multiple  $\Gamma$  matrices, where each one corresponds to some subset of  $X$ , obtained by partitioning  $X$  based on the values of some attribute  $G$  (e.g. class, cluster id). (2) Assuming dimensions are independent, resulting in a diagonal  $\Gamma$  matrix with  $O(d)$  non-zero entries. Both generalizations can be combined resulting in multiple data summaries assuming independence or dependence, covering a wide spectrum of models. Generalization (1) can potentially benefit classification models like Naive Bayes and Linear

Discriminant Analysis Generalization (2) is the gateway to compute simple descriptive statistics like the mean and standard deviation on multiple data subsets and also clustering models. In general, both generalizations require iterative methods, reading the data set at every iteration. Therefore, their solution is different from the one-pass solution presented in this work and it is an issue for future work.

We now present customized algorithms to compute each model  $\Theta$  under our unifying 2-phase method.

### 4.3.1 Principal Component Analysis (PCA)

The objective of PCA is to reduce the noise and redundancy of dimensions by re-expressing the data set  $X$  on a new orthogonal basis, which is a linear combination of the original dimensions basis. In this case  $X$  has  $d$  potentially correlated dimensions but no  $Y$ . In general, PCA is computed on the covariance or the correlation matrix of the data set [14].

We focus on computing PCA on the covariance matrix  $V$  or the correlation matrix  $\rho$  [14] of  $X$ , both symmetrical matrices. The PCA model  $\Theta$  consists of  $U$ , a set of  $d$  orthogonal vectors, called the principal components of the data set, ordered in decreasing order by their length (variance) and a diagonal matrix  $D^2$ , with the squared eigenvalues (lengths). PCA is computed with an eigen decomposition of  $V$  or  $\rho$ , whose solution is the factorization  $\rho = UD^2U^T = (UD^2U^T)^T$ . We stress we are not computing the SVD factorization  $X = UDV^T$ , where  $U \neq V$  as that is a more general numerical analysis problem [7, 14]. The correlation matrix  $\rho$  is generally preferred because it normalizes dimensions. Notice that when  $X$  has been normalized

(i.e. with a z-score, centering at the origin subtracting  $\mu$  and rescaling variables by standard deviation  $\sigma_j$ )  $\rho = XX^T/n$ . Therefore, given  $\rho$ , its eigen decomposition [14] is a symmetric matrix factorization expressed as  $\rho = (XX^T)/n = UD^2U^T$ . In general, only  $k$  principal components ( $k < d$ ) carry to most important information about variance of the data set and the remaining  $d - k$  components can be safely discarded to reduce  $d$ . The actual reduction of  $d$  to some lower dimensionality  $k$  (i.e., the  $k$  principal components) requires an additional matrix multiplication, which is simpler and faster than computing  $\Theta$ .

As mentioned in Chapter earlier, PCA can be computed solving SVD on the correlation matrix. To compute PCA we rewrite the correlation matrix equation based on the sufficient statistics in  $\Gamma$ :

$$\rho_{ab} = (nQ_{ab} - L_a L_b) / (\sqrt{nQ_{aa} - L_a^2} \sqrt{nQ_{bb} - L_b^2}) \quad (4.6)$$

The PCA algorithm two phases are:

1. Compute  $\Gamma$ .
2. Compute  $\rho$  (or covariance matrix  $V$ ), solve SVD of  $\rho$  (or  $V$ ). Optional: select the  $k$  principal components (generally whose eigen-value is at least 1).

We omit discussion of the actual dimensionality reduction of  $X$  to a  $k$ -dimensional data set. This computation is straightforward requiring only a matrix multiplication between a  $d \times k$  matrix  $\Gamma$  derived from SVD and  $X$ .

### 4.3.2 Linear Regression (LR)

Let  $X = \{x_1, \dots, x_n\}$  be a set of  $n$  data points with  $d$  explanatory variables and  $Y = \{y_1, \dots, y_n\}$  be a set of values such that each  $x_i$  is associated to  $y_i$ . The linear regression model characterizes a linear relationship between the dependent variable and the  $d$  explanatory variables.

Using matrix notation, the data set is represented by matrices  $Y$  ( $1 \times n$ ) and  $X$  ( $d \times n$ ), and the standard definition of a linear regression model is:

$$Y = \beta^T \mathbf{X} + \epsilon \quad (4.7)$$

where  $\beta = [\beta_0, \dots, \beta_d]$  is the column-vector of regression coefficients,  $\epsilon$  represents Gaussian error and for mathematical convenience  $\mathbf{X}$  represents  $X$  augmented with an extra row of  $n$  1s stored on dimension  $X_0$ . The vector  $\beta$  is usually estimated using the ordinary least squares method [19], whose solution is:

$$\hat{\beta} = (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}Y^T \quad (4.8)$$

As introduced above,  $\Gamma$  contains those 2 partial matrix products and so  $\hat{\beta} = Q^{-1}(\mathbf{X}Y^T)$ . Therefore, the LR algorithm becomes:

1. Compute  $\Gamma$ .
2. Solve  $\hat{\beta}$  exploiting  $\Gamma$ .

### 4.3.3 Variable Selection (VS)

We now focus on one of the hardest analytical problems: variable selection in statistics, feature selection in machine learning. Because we have previously introduced linear regression, we will study variable selection in this model. However, our algorithms have the potential to be used in other models such as probit models [14], Bayesian classification [14] and logistic regression [14], all of which have similar intermediate matrix computations.

The input is  $X$ , like LR, but it is necessary to add model parameters. In general, a linear regression model is difficult to interpret when  $d$  is large, say more than 20 variables. The search for the best subsets of explanatory variables that are good predictors of  $Y$  is called variable selection. The assumption of this search is that the data set contains variables that are redundant, or that they have low predictive accuracy and thus they can be safely excluded from the model. When  $d$  is high, an exhaustive search on the  $2^d$  subsets of variables is computationally intractable for even a moderately large  $d$ . Traditional greedy methods, for instance, stepwise selection [14] can only find one subset of variables, likely far from the optimal, but the probabilistic Bayesian approach aims to identify several promising subsets of variables, which in theory can be closer to the optimal. In this work, the set of selected variables will be represented by a  $d$ -dimensional vector  $\gamma \in \{0, 1\}^d$ , such that  $\gamma_j = 1$  if the variable  $j$  is selected and  $\gamma_j = 0$  otherwise. We denote by  $\Theta_\gamma$  the model that selects  $k$  out of the  $d$  variables, corresponding to the vector  $\gamma$ . The dot product  $k = \gamma^T \cdot \gamma$  indicates how many variables were selected. Throughout this paper we will use  $\gamma$  as an index to project matrices on selected variables such as  $\beta_\gamma$

and  $\mathbf{X}_\gamma$ .

## MCMC Methods: the Gibbs Sampler for VS

We use the Bayesian formulation and Gibbs sampler introduced in [19] based on a Zellner G-prior, which has the outstanding feature of sharing several intermediate matrix computations with the other models defined above. Given a model  $\Theta$  with  $K$  parameters, one of the parameters is sampled from its prior distribution, while the others  $K-1$  parameters remain fixed. The set of variables selected at iteration  $i$  of the sequence of observations is denoted by a vector  $\gamma^{[i]}$ . The Gibbs sampler for variable selection generates the Markov chain:  $\gamma^{[0]}, \dots, \gamma^{[N]}$ , in which it is expected the best variable subsets will appear more frequently. At each iteration  $I$ , the Gibbs sampler generates  $\gamma^{[I]}$  based on the previous state of the Markov chain,  $\gamma^{[I-1]}$ . Since each coordinate  $\gamma_j^{[I-1]}$  is a parameter of model  $\Theta_\gamma$ , the Gibbs sampler visits all  $d$  entries of  $\gamma^{[I]}$ . In more detail, for every variable  $\gamma_j^{[i]}$  the normalized probabilities  $p(\gamma_j^{[i]} = 0)$  and  $p(\gamma_j^{[i]} = 1)$  are calculated. Based on these probabilities either  $\gamma_j^{[i]} = 0$  or  $\gamma_j^{[i]} = 1$  is chosen by sampling and the  $j$ th position of vector  $\gamma_j^{[i]}$  is updated accordingly. The Markov chain becomes stable after an initial number of iterations, known as the *burn-in* period  $B$ . Then the MCMC method iterates for a large number of iterations until the posterior probabilities of  $\gamma$  have been sufficiently explored (essentially a CPU bound computation, which requires many trial and error runs).

Recall we perform variable selection with an MCMC method, potentially requiring thousands of iterations. We base the computation of  $\pi(\gamma|X, Y)$  on the Zellner G-prior due to [19]. Zellner's G-prior relies on a conditional Gaussian prior for  $\beta$  and an improper (Jeffrey's) prior for  $\sigma^2$  [19], with parameters  $c$  and  $\tilde{\beta}$ . The most

common prior probability for  $\gamma$  is the uniform prior  $\pi(\gamma|X) = 2^{-d}$  [19]. Under Zellner's G-prior, the parameter  $c$  influences the number of variables selected by the model: large  $c$  values lead to parsimonious models (few variables), whereas small values of  $c$  promote saturated models (many variables). On the other hand,  $\tilde{\beta}$  is a hyper-parameter [19] (i.e. a new parameter not belonging to the original linear regression model) to impose a Gaussian on  $\beta$ .

$$\beta|\sigma^2, X \sim \mathcal{N}_{k+1} \left( \tilde{\beta}, c\sigma^2 (XX^T)^{-1} \right) \quad (4.9)$$

$$\sigma^2 \sim \pi(\sigma^2|X) \propto \sigma^{-2} \quad (4.10)$$

For notational convenience  $\mathbf{X}_\gamma, \tilde{\beta}_\gamma, \mathbf{Q}_\gamma$  mean projecting  $\mathbf{X}, \tilde{\beta}, \mathbf{Q}$  on selected variables. Recall  $\mathbf{X} = [1, X]$ . The full equation to get the posteriors of  $\gamma$  involves:

$$\begin{aligned} \pi(\gamma|X, Y) \propto & (c+1)^{-\frac{k+1}{2}} (YY^T \\ & - \frac{c}{c+1} (Y\mathbf{X}_\gamma^T)(\mathbf{Q}_\gamma)^{-1}(\mathbf{X}_\gamma Y^T) \\ & - (c+1)^{-1} \tilde{\beta}_\gamma^T \mathbf{Q}_\gamma \tilde{\beta}_\gamma)^{-n/2} \end{aligned}$$

The Zellner G-prior is computationally convenient for big data, because equations can be expressed in terms of  $\Gamma$ . We will use the sufficient statistics by equations 4.1, 4.2 and 4.3 in order to avoid recomputing those matrix products at each iteration. Therefore, our optimized algorithm becomes

1. Compute  $\Gamma$ , which contains  $\mathbf{Q}_\gamma, \mathbf{X}_\gamma Y^T$  and  $YY^T$ .



2. Iterate the Gibbs sampler a sufficiently large number of iterations to explore  $\pi(\gamma|X, Y)$ .

# Chapter 5

## Data Summarization in Parallel Database Systems

In this chapter we elaborate on the system implementation details in two parallel database systems with  $N$  processing units in a shared-nothing architecture (i.e.,  $N$  nodes, each with its own memory and disk space). We also discuss in detail the general and system-dependent optimizations.

### 5.1 Data summarization in SciDB

A parallel array DBMS enables the manipulation of multidimensional arrays, removing RAM limitations in a large matrix computation. The fundamental difference between SciDB and a row or column DBMS is storage by chunk instead of block. Such chunk storage requires different programming and optimization, compared to SQL queries or UDFs. Moreover, it is feasible to develop an operator that receives a matrix as input and produces a matrix as output, something not easy to do in SQL.

SciDB organizes array storage by chunks [28], where each chunk is a large multidimensional block on secondary storage that stores a preferably large fraction of cells (values) for one numeric attribute (i.e., a subarray). Array storage in SciDB can be visualized as a grid of rectangular chunks (tiles in 2D), where each chunk comprises a rectangular area of contiguous cells. Storage for dense and sparse arrays

is different, but SciDB allows an efficient migration to a dense representation as a matrix gets denser.

### 5.1.1 Programming Mechanisms in SciDB

The major question is how to program the computation of  $\Gamma$  in SciDB. Choices are: (1) exporting (with a bulk mechanism) the  $X$  matrix to a linear algebra package (LAPACK); (2) using relational style queries on arrays in the AQL (Array Query Language) language provided by the DBMS; (3) composing existing array operators in the AFL (array functional language) language; (4) developing a new array operator, tailored to  $\Gamma$ . Exporting the input matrix is a bad idea, as it defeats the purpose of the DBMS, no matter how fast the computation can happen outside. The query language is a natural second option, but it is inefficient as it requires a similar evaluation query plan as an SQL query with a costly relational join operation. In fact, we will show such join is slow even on a fast column DBMS. Existing matrix operators are inefficient and limited by RAM as  $\Gamma$  is a matrix product that is difficult to optimize because the input is a large rectangular matrix. These alternatives were experimentally compared in [24] and our new matrix summarization operator was the winner.

### 5.1.2 The Gamma Array Operator to Summarize a Matrix

As pointed out by previous research [23], it is easy to compute only  $n, L$ , without  $Q$ , because they only require sums of a value (counting) and incrementally summing a  $d$ -vector. That is,  $Q$  makes the computation more demanding. Recalling that

$\Gamma$  contains  $n, L, Q$  as submatrices, the main challenge is to compute this matrix product:

$$\Gamma = \mathbf{Z}\mathbf{Z}^T = \sum_{i=1}^n z_i \cdot z_i^T,$$

which requires cross-products between all dimension pairs, an  $O(d^2n)$  computation. That is, we study the summarization operator purely as a matrix multiplication that subsumes all previous approaches. A fundamental aspect is to optimize computation when  $\mathbf{X}$  (and therefore  $\mathbf{Z}$ ) is sparse; any multiplication by zero returns zero. Therefore, when computing  $\mathbf{Z}\mathbf{Z}^T$  a multiplication should be evaluated only when both vector  $z_i$  entries are different from zero. Since dense and sparse matrices have different storage and different processing in main memory this leads specialized summarization operators for each of them.

### Dense Matrix Operator

The dense operator requires that all data of  $z_i$  ( $x_i$  in consequence) fits in one chunk. In other words,  $z_i$  cannot be partitioned into several chunks in spanned storage. This requirement is important to accelerate I/O: larger chunks improve I/O because our  $\Gamma$  computation requires a full scan on the data set, but no join algorithms. Otherwise, if  $d$  is so large that it is necessary to store  $x_i$  spanning several chunks, then it would be necessary to call a join algorithm, joining  $X$  with itself.

Smaller chunks would trigger seeks, resulting in worse performance. Fitting many points  $x_i$  in one chunk is not an unreasonable assumption because SciDB favors fairly large chunk sizes, typically ranging above 8MB. Another major requirement

and advantage is that the operator must work in parallel. It is a requirement that complicates programming because it is essential to design and develop the operator in such a way that  $X$  can be evenly partitioned across  $N$  processing nodes, with minimal or no need of synchronization overhead. The advantage is, of course, the ability to scale processing to larger  $N$  as  $n$  grows. Our  $\Gamma$  operator works fully in parallel with a partition of  $X$  into  $N$  subsets  $X^{[1]} \cup X^{[2]} \cup \dots \cup X^{[N]} = X$ , where we independently compute  $\Gamma^{[I]}$  on  $X^{[I]}$  for each node  $I$  (we use this notation to avoid confusion with matrix powers and matrix entries). In SciDB terms, each worker  $I$  will compute  $\Gamma^{[I]}$ . When all workers are done the coordinator node will gather all results and compute a global  $\Gamma = \Gamma^{[1]} + \dots + \Gamma^{[N]}$  with  $O(d^2)$  communication overhead per node (much smaller than  $O(dn)$ ). This is essentially a natural parallel computation [27], coming from the fact that we can push the actual multiplication of  $z_i \cdot z_i^T$  into the array operator. Given the additive properties of  $\Gamma$  the same algorithm is applied on each node  $I = 1 \dots N$  to get  $\Gamma^{[I]}$ , combining all partial results  $\Gamma = \sum_I \Gamma^{[I]}$  in the coordinator node. Our parallel algorithm to compute  $\Gamma$  is below.

### **Sparse Matrix Operator**

A sparse matrix uses less space on secondary storage, resulting in faster I/O and a smaller footprint in RAM per point. A second advantage is that since more cells fit in one chunk  $d$  can be higher compared to the dense operator. Therefore, assuming a density threshold  $\psi$ ,  $x_i$  can have higher dimensionality being able to fit in one chunk. If  $p$  is the maximum dimensionality for a dense matrix representation then  $d = \lceil p/\psi \rceil$  can be the maximum dimensionality for a sparse matrix. For instance, if  $\psi = 0.01$  then  $d = 100 \times p$ .

```

/* X on secondary storage,  $\Gamma$  in main memory */
Data:  $X = \{x_1, x_2, \dots, x_n\} = X^{[1]} \cup \dots \cup X^{[N]}$ 
Output:  $\Gamma$ 

/* Parallel for:
   N nodes working in parallel on partition  $X^{[I]}$  */
for  $I = 1 \dots N$  do
   $\Gamma^{[I]} \leftarrow \mathbf{0}$ 
  for  $x_i \in X^{[I]}$  do
     $z_i = [1, x_i, y_i]$ 
    for  $a = 0 \dots d + 1$  do
      for  $b = 0 \dots a$  do
         $\Gamma_{ab}^{[I]} \leftarrow \Gamma_{ab}^{[I]} + z_{ia} * z_{ib}$ 
      end
    end
  end
end

/* send local summaries to coordinator node */
 $\Gamma = \sum_{I=1}^N \Gamma^{[I]}$ 

```

**Algorithm 1:** Parallel dense algorithm to compute lower triangular  $\Gamma$  on a dense matrix.

In a similar manner to a dense matrix, we assume  $x_i$  fits in one chunk. Moreover, since  $x_i$  is a sparse vector its coordinates reside on the same chunk improving I/O locality. Under this assumption, the algorithm for a sparse matrix has the same outer loop on  $i$  to scan partition  $X^{[i]}$ , compared to the algorithm for a dense matrix. The first major difference is that we *dynamically* build a main memory sparse vector representation for  $z_i$ , following a similar scheme to LAPACK [8], whose number of non-zero entries is  $k = |z_i|$ . The loop to build  $z_i$  is  $O(d)$  (and not  $O(k)$ ) because SciDB does not store matrix entries as (subscript,value) pairs. Instead, SciDB uses Run-Length Encoding for a sequence of zeroes, which consumes more disk space, but allows an easy gradual change from sparse to dense representation. Then the second major difference is the two inner loops, which go  $a = 1 \dots k, b = 1 \dots k$  which result in  $O(k^2)$  flops. Such optimization results in a much better time complexity  $O(k^2n) < O(d^2n)$ . Assume that  $X$  is hyper-sparse [6] so that  $k^2 \leq d$  (i.e.  $k = O(\sqrt{d})$ ; the ratio of number of non-zero entries to  $d$  asymptotically approaches zero). Then time complexity becomes  $O(dn)$ , which is significantly lower than  $O(d^2n)$ . In a similar manner to the dense operator, there is  $O(d^2)$  communication overhead per node (much smaller than  $O(dn)$  if the data set had to be transferred) to get the global summary matrix  $\Gamma$  in time  $O(d^2N)$  using a locking mechanism to update the matrix. Finally,  $n_{nz} \leq \sqrt{dn}$  can be used as a threshold to use the sparse algorithm, where  $n_{nz}$  is the number of non-zero entries in  $X$ .

The sparse operator brings two computational aspects into question. The first issue is being able to maintain  $\Gamma$  in RAM for high  $d$ . The second one is developing yet a third operator for a sparse  $\Gamma$  matrix. We emphasize again that for most big data analytics problems  $d < n$ . For the first issue it is reasonable to maintain  $\Gamma$  for fairly

```

/* X on secondary storage,  $\Gamma$  in main memory */
Data:  $X = \{x_1, x_2, \dots, x_n\} = X^{[1]} \cup \dots \cup X^{[N]}$ 
Output:  $\Gamma$ 

/* Parallel for:
   N nodes working in parallel on partition  $X^{[I]}$  */
for  $I = 1 \dots N$  do
   $\Gamma^{[I]} \leftarrow \mathbf{0}$ 
  for  $i = 1 \dots |X^{[I]}|$  do
     $k \leftarrow 1$  /*  $z_{i0} = 1$  */
     $z_i \leftarrow [1, x_i, y_i]$ 
    for  $j = 0 \dots d + 1$  do
      if  $z_{ij} \neq 0$  then
         $v_k \leftarrow z_{ij}$ 
         $i_k \leftarrow j$ 
         $k \leftarrow k + 1$ 
      end
    end
    /*  $k = |z_i|$  */
    for  $a = 1 \dots k$  do
      for  $b = 1 \dots a$  do
         $\Gamma_{i_a i_b}^{[I]} \leftarrow \Gamma_{i_a i_b}^{[I]} + v_{i_a} * v_{i_b}$ 
      end
    end
  end
end

/* send local summaries to coordinator node */
 $\Gamma = \sum_{I=1}^N \Gamma^{[I]}$ 

```

**Algorithm 2:** Parallel sparse algorithm to compute lower triangular  $\Gamma$  on a sparse matrix.



high  $d$  in RAM given the rapid growth in RAM volume (we experimentally show this is feasible for  $d \approx 1000$ ). For the second issue, in general, a sparse representation for  $\Gamma$  does not make sense because that would imply the data set contains zero-constant dimensions to start with (i.e., useless). Moreover, any constant dimension should be detected and eliminated, alerting the user. For such cases it makes sense to discard constant dimensions on a pre-processing pass computing only  $n, L$  and the diagonal of  $Q$ . That is, computing only row 0 (column 0) and the diagonal of  $\Gamma$  (no dimension cross-products). Then the user can safely discard any dimension  $j$  whose  $\sigma_j = 0$ . This a consequence of  $\Gamma$  containing an extended set of sufficient statistics from which we can derive  $\mu_j, \sigma_j$  for dimension  $j$ .

### 5.1.3 Additional Considerations

#### Data quality: Constant dimensions and missing values

Our  $\Gamma$  summarization matrix can help detecting and fixing data quality problems. Based on the fact that  $\mu_a, \sigma_a$  can be easily derived from  $\Gamma$ , the  $\Gamma$  matrix can be used to eliminate constant and near-constant dimensions where  $\sigma_a \approx 0$ .

Consider the trivial, but common, case that  $\Gamma_{aa} = 0$  for  $a = 1 \dots d+1$  (remember that  $\Gamma_{00} = n$ ). This means that dimension  $a$  is completely full of zeroes: clearly such dimension is useless and should be eliminated from further consideration. In a more general case, consider constant dimensions with non-zero values. If dimension  $a$  is constant then the variance of dimension  $a$  is zero, then all covariances  $\Sigma_{ab}$  or correlations  $\rho_{ab}$  are undefined (null in SQL). The reason is straightforward: a constant dimension provides no statistical insight to the analyst (other than alerting

the analyst that something may be wrong in data collection). Moreover, PCA and linear regression require eliminating constant dimensions because their intermediate matrices cannot be factorized or inverted. In short, we assume that before computing model  $\Theta$  all constant and near-constant dimensions have been eliminated, and any missing values have been replaced or their points eliminated.

## Systems Aspects

Our operators were programmed in C++ with the API provided by SciDB. The operator was programmed as an array aggregation (Figure 5.1), which has similarities with an aggregate UDF in a relational DBMS [23]. The similarities are: the ability to update variables in main memory, seamless integration with a scan over the input array, parallel execution with just one synchronization barrier at the end, being programmed in an efficient language (the same as the DBMS: C++), similar processing phases (initialization, increment, final and return), and not being able to call other array operators within the operator source code. However, as fundamental differences we must mention: our array operator receives array chunks as input (not rows) and it produces an array as output, not a scalar value. Therefore, it is unnecessary to do any data type conversion at the end as it happens in a relational DBMS [23]. The ability to define a parallel array operator that receives an array as input and produces an array as output was fundamental for us.

We now discuss why the dense and sparse operators working together with R can compute models faster than alternative approaches. Specifically, our array operator is much faster than an implementation on Hadoop-MapReduce systems (e.g. Spark [34]) and it is much faster than any SQL-based engine, including UDFs. This is a

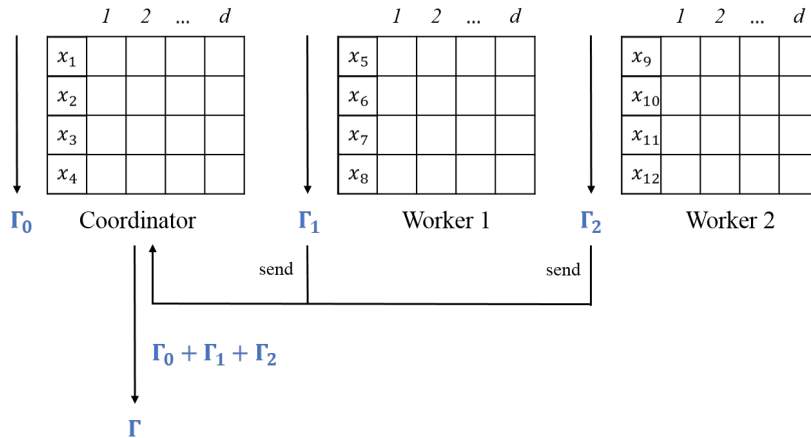


Figure 5.1: How the Gamma operator works in SciDB

list of main reasons. Our operator works in a single array scan with large disk blocks (chunks) and these chunks are fairly large ( $>8\text{MB}$ ) [28] (and growing in the future as RAM grows). The operators work with compiled C++ code: there is no overhead or delay parsing a query, query optimization at run-time or ODBC communication overhead. Moreover, compiled C++ code is faster than Java byte code and faster and more memory-efficient than Java compiled code (widely used on Hadoop). Since we use an array operator that takes an array as input and an array as output there is no overhead converting matrices between different representations. Our operator cannot be as efficient as compiled FORTRAN code used by LAPACK for small matrices, where small means a matrix that can fit in cache memory in the CPU. However, our operator scales much better as matrices get bigger and they surpass cache memory size. In fact, we will show our operator surpasses LAPACK speed for very large

sparse matrices, which is the state of the art for high performance computing. Such comparison settles any speed concerns about our operator. On the parallel processing side, the input matrix storage layout is key to enable a fully parallel computation without needing to send data from one node (or processing thread) to another one. Our main assumption is that  $x_i$  can fit completely in one chunk: we will show experiments where  $d = 800$ , a truly high dimensionality. Moreover, we can solve even higher  $d$  problems with sparse matrices, the norm in big data analytics. While it is true, data shuffling in a distributed system may be needed for high dimensional data we will experimentally show we can solve fairly high  $d$  problems. Considering pure C++ code reading a binary file our operator will likely be faster simply because SciDB has an efficient I/O subsystems, processing is multi-threaded and access is block-based. In other words, a well-crafted C++ program would need to mimic or redo the basic SciDB array storage. Therefore, our system runs much faster and is more scalable than any other C++, LAPACK, SQL or Hadoop-based system. Our experiments will prove that is the case.

## 5.2 Data Summarization in Columnar DBMS

We start by discussing storage in the columnar DBMS, where indexing is not required. The data set and model matrices are stored in tables, having a primary key with a specific combination of matrix subscripts, depending on their storage layout. We use two fundamental data set layouts for  $X$ : horizontal and vertical. The first layout is a standard tabular representation which has  $n$  rows, each having  $d$  columns. The second layout is based on a pivoted table with one attribute value

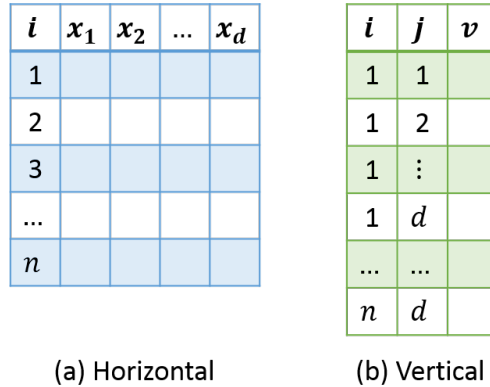


Figure 5.2: The horizontal and vertical layout for  $X$  in the columnar DBMS.

per row and up to  $dn$  rows, which enables efficient processing of sparse matrices and which also removes DBMS limitations on the maximum number of columns. The I/O efficiency of each layout depends on matrix sparsity and the specific query physical operator (scan, join, sort). Figure 5.2 illustrates the two types of layout.

In order to achieve better performance, the data set will be sorted by the matrix subscripts, and will be segmented and distributed to all the  $N$  processing nodes evenly according to the value of a hash function on the row number  $i$ . This process will be automatic in the columnar DBMS during the data loading phase once we properly define the table at creation time. The following SQL command shows a proper way to define the input table  $X$  in Vertica:

```
CREATE TABLE X (
     $i$  INT NOT NULL ENCODING RLE,
     $j$  INT NOT NULL ENCODING COMMONDELTA.COMP,
```

```
v FLOAT)  
ORDER BY i, j  
SEGMENTED BY MODULARHASH(i) ALL NODES;
```

Computing matrix equations with SQL queries has proven to be hard to optimize [21]. To get started, it is not possible to create fixed queries that can efficiently evaluate mathematical equations because the table definition for the input data set is not fixed (especially for the horizontal layout with  $d$  columns). Second, evaluating matrix equations with pure SQL queries usually involves table joins, which are always quite expensive for DBMSs especially when they are evaluated in parallel. Therefore, SQL code generation at run-time from R and developing specialized mathematical UDFs to bypass table joins are necessary. Statistical methods and data mining algorithms are also programmed with R in order to take full advantage of its rich library of models and techniques.

The R program is based on two sets of parameters. The first set of parameters controls SQL code generation (i.e. table name and layout etc.) and database systems optimizations, which are tailored to each algorithm (with some essential math optimizations across multiple models). The second set of parameters controls the mathematical behavior of each technique. Examples include the number of clusters, tolerance threshold for convergence, numeric stability issues, and so on.

We must emphasize that SQL code generated by R represents a Turing-complete language. That is, we can evaluate any equation with matrices, no matter how complex it is.

As we have discussed in Chapter 4.1,  $\Gamma = Z \cdot Z^T$  [24].  $\Gamma$  can be computed by a single SQL query using self-joins like the following:

```

SELECT
  a.j AS i,
  b.j AS j,
  SUM(a.v * b.v) AS v
FROM X a JOIN X b ON a.i = b.i
GROUP BY a.j, b.j;

```

To bypass the expensive self-join in the SQL query, we can also compute the  $\Gamma$  by calling a user-defined transform function that builds vector  $z_i$  and updates  $\Gamma$  in parallel and incrementally:

```

SELECT Gamma(i, j, v USING PARAMETERS d = d)
OVER (PARTITION BY MOD(i, N)
ORDER BY i, j)
FROM X;

```

### 5.3 Time Complexity and Parallel Speedup

Our analysis is based on the 2-phase algorithm introduced in Section 4.2. We start with time complexity. Phase 1, which corresponds to the summarization matrix operator, is the most important. Computing  $\Gamma$  with the dense matrix operator is  $O(d^2n)$ . On the other hand, computing  $\Gamma$  with the sparse matrix operator is  $O(k^2n)$  for the average case assuming  $k$  entries from  $x_i$  are non-zero. Assuming  $X$  is hyper-sparse  $k^2 = O(d)$  then the matrix operator is  $O(dn)$  on average. Space required by  $\Gamma$  in main memory with a dense representation is  $O(d^2)$ . In Phase 2 we take advantage of  $\Gamma$  to accelerate computations involving  $X$ . Because we are computing matrix factorizations derived from  $\Gamma$  time is  $\Omega(d^3)$ , which for a dense matrix it may

approach  $O(d^4)$ , when the number of iterations in the factorization numerical method is proportional to  $d$ . In short, time complexity for Phase 2 for the models we consider does not depend on  $n$ . I/O cost is just reading  $X$  with the corresponding number of chunk I/Os in time  $O(dn)$  for dense matrix storage and  $O(kn)$  for sparse matrix storage.

We now analyze parallel speedup assuming  $N$  processors on a shared-nothing parallel architecture. Recall we assume  $d < n$  and  $n \rightarrow \infty$ . For Phase 1 each chunk of  $X$  is hashed to  $N$  processors. Because we assume  $x_i$  fits in one chunk this results in  $x_i$  being hashed as well. Then our matrix operator can compute each outer product  $x_i \cdot x_i^T$  in parallel, resulting in  $O(d^2n/N)$  work per processor for a dense matrix and  $O(dn/N)$  work per processor for a sparse matrix. Both dense and sparse matrix operators become optimal as  $N$  increases or  $n \rightarrow \infty$ . We now consider Phase 2, which involves numerical methods from linear algebra. Since  $\Gamma$  significantly reduces problem complexity and the numerical methods used to solve SVD, least squares and matrix inversion are mathematically complex we assume matrix factorizations are computed on one node ( $N = 1$ ) calling the efficient LAPACK library. In other words, it is not worth reprogramming SVD, least squares or matrix inversion to run in parallel across  $N$  nodes, but they can indeed run in parallel on one node with a multicore CPU, taking advantage of LAPACK. Notice the MCMC method used by VS is difficult to parallelize across iterations since each iteration depends on the previous one, but it makes sense computing each iteration, requiring a matrix inversion, as fast as possible. Therefore, for Phase 2 there are two choices for computation: totally sequential or parallel with scale up (increasing number of cores and threads). In short, for Phase 2 we rely on the parallel and numeric capabilities provided by



LAPACK, which is a gold standard for numeric accuracy and which provides superb performance on one node with multicore CPUs.

# Chapter 6

## Experimental Evaluation

We present experiments focusing on scalability and parallel processing. Because our operator and the 2-phase algorithm do not change the accuracy of matrices and final results it is unnecessary to measure statistical or numeric accuracy. We should emphasize that our summarization operator was used to produce the  $\Gamma$  matrix to be consumed by the R package. That is, the final model  $\Theta$  computation happened in R. Therefore, the model computed by R alone and the model computed by SciDB in tandem with R are the same.

This is an overview of experiments with a focus on scalability and parallel processing. We first review existing system and justify our choice of specific systems for comparison. We start our evaluation presenting time benchmarking experiments comparing our operator with popular analytic systems, including the R package and a fast columnar DBMS on a single computer with a powerful multicore CPU. We should stress R and many math packages internally use LAPACK for the most numerically intensive computations. Therefore, in order to test our matrix operators with large matrices, we also make a careful comparison with the LAPACK [8, 7] library (the state of the art in numerical linear algebra) using its Intel MKL variant (the fastest variant for Intel multicore CPUs). Finally, we analyze scalability and speedup on a large cluster in the cloud with large data sets comparing our R+SciDB system with Spark, currently the most popular analytic platform from the Hadoop stack.

Table 6.1: Data sets.

Data set	$d$	$n$	Description
KDDnet	38	5M	Network intrusion detect (KDD Cup)
Gene	12506	248	Gene data to predict cancer survival

Table 6.2: Hardware and Operating System.

Item	Local	Cloud (Amazon)
OS	Linux Ubuntu	Linux Ubuntu
CPU	4 cores: 2.15 GHz	2 cores (2VCPU): 2GHz/worker 4 cores (4VCPU): 2GHz/coord.
$N$ nodes	1	1, 10, 100
threads	SciDB: 2 instances Spark: NA	SciDB: 2,20,200 instances Spark: 2,20,200 executors
RAM	4 GB	7.5 GB/node (SciDB/Spark) 15 GB/coordinator (Spark)
Storage	3 TB	10 TB

## 6.1 Experimental Setup

We used two data sets: the network intrusion data set and a microarray genomics data set, summarized in Table 6.1, obtained from the KDD Cup repository (KDD 99), and a cancer data set repository, respectively. We sampled and replicated the KDDnet to get varying  $n$  (data set size) and  $d$  (dimensionality), without altering its statistical properties. Its columns were replicated three times and rows were replicated to get  $d = 100, n = 10M$ . Based on this “reference” data set rows were sampled and replicated in log-10 scale and columns were replicated in log-2 scale. On the other hand, the Gene data set was left as is.

### 6.1.1 Mathematical Models Functions and Parameters

For PCA and LR we used the R default parameter settings to get accurate results. On the other hand, for VS we used a small number of iterations for the MCMC Gibbs sampler since this algorithm is CPU bound (after being optimized) and we were interested in comparing performance, rather than getting a very accurate Bayesian model, which would require many trial/error experiments and tens of thousands of iterations (i.e., exploring the posterior probability and evaluating convergence of the Markov Chain). For Spark we used the available functions it had for matrix multiplication, LR and PCA and MLlib, the best library currently available. We should mention Spark did not offer functions for variable selection. For LR and PCA we used the recommended number of iterations by Spark.

### 6.1.2 System Setup and Tuning

We tested our operator and competing systems on two hardware configurations (summarized in Table 6.2): (1) one node ( $N = 1$ ) (4-core CPU) and (2)  $N$  nodes (2-core VCPU). Our default system configurations were two instances for SciDB and two executors per node for Spark (i.e., 2 threads per node). For  $N = 1$  data sets were copied to the Linux file system and then loaded into the respective DBMS (array, column). For the  $N$  nodes we used the Amazon cloud, which already offered optimal configuration for HDFS, and where Spark was pre-installed and tuned. The times to transfer, copy and load the data sets into chunk format for SciDB or the time to copy csv files from Unix to HDFS for Spark are not included in our time measurements, since it is a one-time event. Processing times do include the I/O time to read arrays

in chunk format in SciDB and the time to read the csv files and transform them into RDD format [34], the distributed memory format in Spark.

For the array DBMS and the column DBMS we cleared the buffers (DBMS cache) before each run to make sure measurements include the time to read from disk. We ran each program three times on each system and we report the average. When the R system crashed, mainly due to RAM limitations, we report “fail”. For Spark we made sure data set was cleared from distributed memory before each run, but Spark could cache the data set it for iterative algorithms (LR, PCA). That is, we attempted to make a fair comparison, giving each system the best opportunity. In general, when the computation could not finish in 30 minutes it was automatically stopped and we report “stop”. We show execution times in seconds. In some isolated cases we let systems run for 1 hour.

## 6.2 Comparing Analytic Systems ( $N = 1$ node)

### 6.2.1 Comparing R+SciDB and R alone

We start by comparing our proposed hybrid solution combining the array DBMS and the R package with the R package alone, assuming a dense matrix as input. We do not compare with a sparse matrix storage in R because most statistical models in R are computed with an input matrix with a dense layout. Moreover, a sparse data set requires a significantly different storage structure in RAM, which would require making at least two passes to convert the data set from dense to sparse storage. For PCA we called the `princomp()` function and for LR we called `solve()`, both standard

R function calls widely used by R analysts. For VS in R we use the R script available from [19] to solve Bayesian variable selection with the same method. The original R script from [19] was so slow that it could not run in under one hour even for  $n=1k$  (i.e., when the data set fits completely in RAM). Therefore, we were forced to incorporate our optimization based on  $\Gamma$ , but as a matrix multiplication. Incorporating a faster optimization based on vector outer products would require developing a C++ function working by block in R, converting the R matrix representation to native C++ arrays, back and forth, which we consider a topic for future research. The SciDB dense matrix operator computes  $\Gamma$ , which is passed to an R script further optimized by us. We would like to mention our time measurements on SciDB and the column DBMS do not include the time to load the data set and apply an array redimension operator to convert it from table to matrix because, in practice, they are a one-time event.

Table 6.3 compares total time to get  $\Theta$  in R alone and the combination DBMS+R. We can see R scales well, but eventually crashes due to reaching RAM capacity. For all models our  $\Gamma$  operator on R+SciDB is faster than R working alone. These times prove the superiority of our operator and our 2-phase algorithm. A combination of factors make R slower: parsing the input text file, single-threaded processing, but mainly not exploiting our summarization matrix.

We now turn our attention to the Gene data set, with very high  $d$ , analyzed on Table 6.4. The authors of [19] provide an R program that computes a state-of-the-art Bayesian model that was our reference Gibbs sampler algorithm. This R program was so slow that we had to incorporate our  $\Gamma$  summarization matrix computed in R as a matrix product. Since  $d > n$  temporary matrices (on selected variables)

Table 6.3: Comparing computation of model  $\Theta$  using R and DBMS+R; dense matrix operator; data set KDDnet; local server; times in secs.

$d$	$n$	PCA		LR		VS	
		R	R+SciDB	R	R+SciDB	R	R+SciDB
10	100k	0.5	0.6	0.5	0.6	3.4	3.6
10	1M	4.8	1.3	5.6	1.3	7.6	4.7
10	10M	45.2	7.0	50.1	7.1	58.8	9.6
10	100M	fail	64.7	fail	69.8	fail	93.4
100	100k	5.7	2.5	6.3	2.6	34.4	14.8
100	1M	61.2	16.8	60.5	16.9	113.0	29.1
100	10M	fail	194.9	fail	194.9	fail	207.2

Table 6.4: Bayesian variable selection (VS) in gene data set ( $d = 12506, n = 248$ ): Comparing R+SciDB with R at 1000 iterations; times in secs.

$p$	R+SciDB	R unoptimized	R optimized with $\Gamma$
100	17	1139	16
200	43	*	43
400	83	*	85
800	199	*	200

were ill-conditioned and the MCMC method did not converge properly. That is, they could not be inverted or factorized. Convergence with such high  $d$  required hundreds of thousands of iterations, becoming prohibitive. Finally, variables with marginal correlation to  $Y$  had little statistical significance. Therefore, we ran an initial variable pre-selection step computing correlations between each variable and  $Y$ , to obtain a reduced dimensionality  $p$ , in which the MCMC method converged properly. As can be seen in Table 6.4,  $\Gamma$  is essential to accelerate computation and R (optimized) was competitive with R+SciDB because the data set could fit in RAM. We emphasize our optimization had a significant impact on R alone, highlighting its generality.

## 6.2.2 Computing $\Gamma$ on R, Array DBMS and Column DBMS

It is natural to wonder if  $\Gamma$  can be computed by the other systems, beyond the array DBMS. With this question in mind, we investigate how  $\Gamma$  can be computed on each system, making a reasonable effort to develop an efficient program within each system constraints.

For R the data set was read from disk and loaded in RAM and then  $\Gamma$  was obtained with a standard matrix multiplication using `%*%` (i.e. the fastest mechanism available, commonly used by R analysts). For the column DBMS we stored  $X$  on a vertical layout with one entry  $X_{ij}$  per row  $(i, j, v)$ , where  $i = 1 \dots n$  and  $j \in \{1 \dots d\}$ . The query to get  $\Gamma$  was a self-join on  $i$  grouping the two dimension subscripts. We defined the table sorted by  $i$  so that all dimension values were clustered on disk; in this manner the query optimizer could use a merge-join, the fastest join possible in time  $O(n)$ . Figure 6.1 compares the three systems: the upper plot compares at low dimensionality  $d = 10$  and the lower plot compares at a high dimensionality  $d=100$ . For sparse matrix experiments zeroes are deleted, resulting in a smaller SQL table and smaller (compressed) array on disk. As can be seen in the plots, the fastest computing mechanisms are our sparse and dense operators, with the sparse operator being at least twice as fast compared to the dense one, with matrices having the same  $d, n$  sizes. In general, our operator is about an order of magnitude faster than R and about two orders of magnitude faster than the SQL query. Table 6.5 complements Figure 6.1, providing a “drill down” view on the specific numbers, demonstrating our operator is both scalable and uniformly more efficient, showing when R failed due to insufficient RAM and indicating when the SQL DBMS had to be stopped due



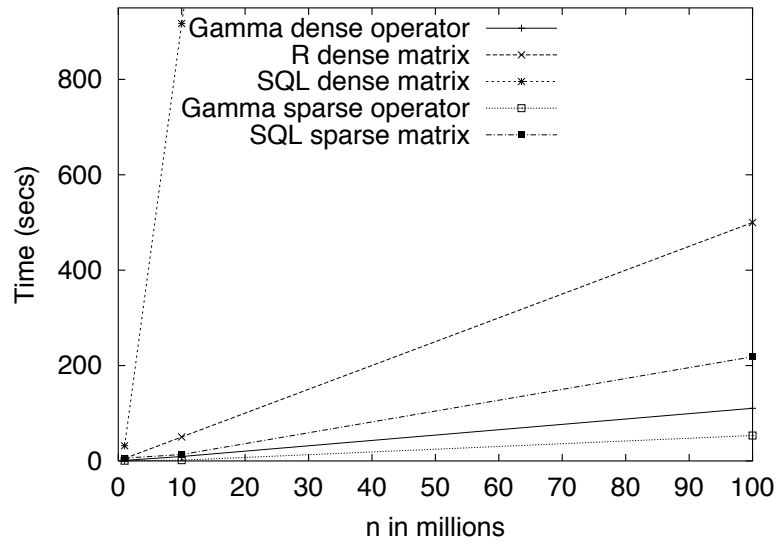
Table 6.5: Comparing matrix summarization algorithms; data set KDDnet; dense ( $dn$  entries) and sparse (zeroes deleted); Local 1/2; times in secs.

$d$	$n$	Dense matrix			Sparse matrix		
		R	SQL	$\Gamma$ dense op.	SQL	$\Gamma$ sparse op.	
					density=12%		
10	100k	0.6	3.8		0.6		0.1
10	1M	5.1	31.9		5.8		0.5
10	10M	50.2	917.2		13.4		1.6
10	100M	fail	stop	110.3	218.4		53.3
					density=27%		
100	100k	6.5	74.0		23.9		1.8
100	1M	44.1	612.5		231.6		10.6
100	10M	fail	stop	332.7	stop		105.2
100	100M	fail	stop	3423.7	stop		1015.3

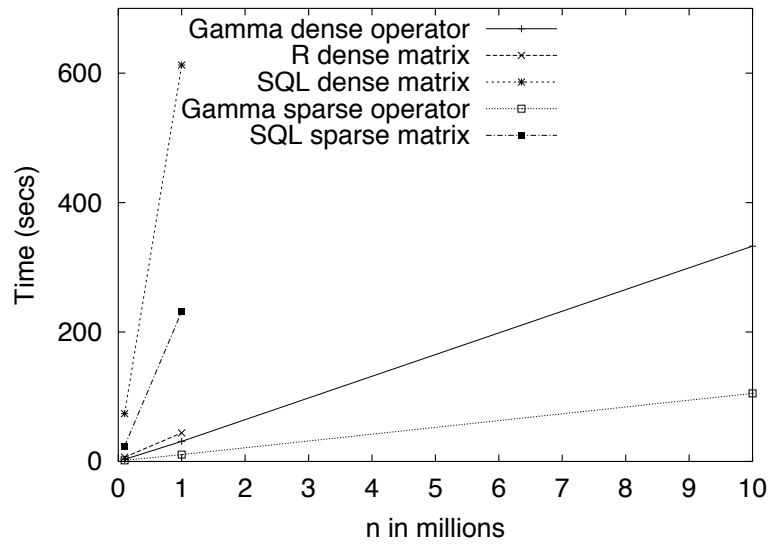
to waiting more than 30 minutes. These results make clear SQL cannot solve dense matrix problems, but it shows promise for sparse matrices of low dimensionality.

### 6.2.3 Comparing $\Gamma$ Algorithm with Fast Matrix Multiplication: SciDB versus Intel MKL (LAPACK)

Here we study in more depth the computation of  $\Gamma$ . As mentioned above, R and many other mathematical packages use the LAPACK library [7, 8], which is both extremely fast and highly accurate. Therefore, it is natural to ask whether it is worth computing  $\Gamma$  with LAPACK (in RAM). Moreover, it is important to understand how LAPACK behaves with truly large matrices, reaching RAM capacity and much larger than L1/L2 cache memory. Parallel LAPACK variants are well tested and used by most systems with dense matrices. On the other hand, as of today developing parallel programs with distributed memory for sparse matrices is an ongoing effort [6]. Here we compare our operator with dense matrix functions available in MKL for multi-core



(a) d=10



(b) d=100

Figure 6.1: Comparing systems to compute  $\Gamma$  on local server: Array DBMS, R, SQL.

Table 6.6: Comparing dense operator (working on disk+RAM), sparse operator (working on disk+RAM) and dense matrix multiplication in MKL BLAS (parallel LAPACK, working only on RAM) to compute  $\Gamma$ ; local server 4GB RAM; synthetic data sets; times in secs.

$n$	density	$d = 100$			$d = 200$			$d = 400$			$d = 800$		
		dense	sparse	MKL	dense	sparse	MKL	dense	sparse	MKL	dense	sparse	MKL
1M	0.1%	27.5	0.2	3.0	94.4	0.2	8.2	374.2	0.4	379.4	1395.4	0.9	fail
1M	1.0%	27.5	0.5	3.0	96.8	1.1	8.2	374.2	3.8	364.7	1404.4	4.0	fail
1M	10.0%	29.4	4.0	3.0	100.7	7.0	8.2	374.2	16.3	367.5	1416.9	46.4	fail
1M	100.0%	30.7	44.0	3.0	106.4	148.8	8.3	374.2	542.3	389.1	1449.4	2159.6	fail
10M	0.1%	290.3	1.1	fail	1003.2	2.0	fail	stop	3.2	fail	stop	stop	fail
10M	1.0%	294.8	6.7	fail	1006.4	12.8	fail	stop	26.6	fail	stop	stop	fail
10M	10.0%	315.0	50.8	fail	1052.3	104.7	fail	stop	195.3	fail	stop	stop	fail
10M	100.0%	323.5	450.6	fail	1058.3	1582.5	fail	stop	stop	fail	stop	stop	fail

CPUs. Later, we compare with a sparse matrix multiplication available in SciDB that works in distributed memory in a parallel cluster (using ScaLAPACK [8]), which requires partitioning the input matrix and using MPI (Message Passing Interface).

Table 6.6 compares our Gamma operators with MKL, the fastest LAPACK parallel version for multicore CPUs. We emphasize MKL is working on RAM only, whereas our operator includes both I/O and CPU time. Therefore, we give MKL an advantage. As a major observation, MKL fails due to insufficient RAM at  $n=1M$ ,  $d=800$  and at  $n=10M$  with all  $d$  values; there is no quick fix to solve this issue other than partitioning  $X$  on secondary storage and developing a new algorithm.

Evidently, MKL is very efficient with small dense matrices and therefore, it is worthless optimizing our operators C++ code for small matrices. The main reason MKL is much faster is because it takes advantage of cache memory (L1 and L2) in the CPU, reading matrix blocks from RAM in the worst case. MKL incorporates sophisticated matrix multiplication algorithms that are very efficient when the matrix is small enough (relative to L1/L2 cache size). Also, MKL incorporates efficient

algorithm to move matrix blocks between cache memory and RAM, back and forth, a low-level optimization out of scope for us. Nevertheless, MKL is marginally faster than our dense matrix operator at  $d=400$ . These results prove our dense matrix operator is slower, but scalable. For sparse matrices the gap gets wider. At the lowest density extreme (density 0.1%) our sparse operator is much faster than MKL and our dense operator. In fact, for large matrices ( $n=1M$ ), density 0.1% and  $d=400$  our sparse operator is *three orders of magnitude* faster than MKL and the dense operator. For  $d=800$  our sparse operator is three orders of magnitude faster than the dense operator. For  $n=10M$  MKL fails with all  $d$  values due to insufficient RAM, the dense operator must be stopped at  $d = 400$ , but our sparse algorithm can still efficiently work at  $d=400$ . Considering that the sparse algorithm is just 50% slower than the dense one with 100% density, it comes out as the overall winner. In summary, our dense and sparse matrix operators working together beat MKL in scalability and speed.

### 6.3 Parallel Processing ( $N = 1, 10, 100$ nodes)

All speedup figures are computed as  $s_N = T_1/T_N$  [8], where  $T_1$  is the time for a purely sequential version and  $T_N$  is the time on  $N$  nodes; ideal speedup means  $S_N \approx N$ .

Table 6.7 presents speedup experiments varying the number of threads on our local server. The goal is to determine the optimal number of threads in the 4-core CPU (i.e., a scale up analysis). As can be seen, as the # of threads grows time decreases, but the time gain shrinks. The optimal number of threads is 2 since those

Table 6.7: Parallel speedup:  $\Gamma$  dense matrix summarization operator with multi-threaded processing varying  $M$ =# SciDB instances ( $N=1$  node, 4-core CPU); times in secs.

$M$	$n=100k$		$n=1M$		$n=100M$	
	time	speedup	time	speedup	time	speedup
1	1.4	1.0	13.0	1.0	146.6	1.0
<b>2</b>	0.7	2.0	6.8	1.9	82.8	1.8
4	0.8	1.8	6.2	2.1	60.8	2.4

Table 6.8: Parallel processing: Comparing R+SciDB and Spark varing  $N$  (# of nodes); data set KDDnet  $d = 38$ ; dense matrix; times in secs.

$N$ nodes	$n$	R+SciDB		Spark					
		Gamma/PCA/LR	speedup	Gramian	speedup	LR (20 iters)	speedup	PCA	speedup
1	1M	7.0	1.0	78	1.0	314	1.0	98	1.0
10	1M	1.0	7.0	84	0.9	246	1.3	101	1.0
100	1M	0.3	25.0	80	1.0	364	0.9	106	0.9
1	10M	69.0	1.0	200	1.0	2524	1.0	352	1.0
10	10M	6.8	10.1	94	2.1	376	6.7	120	2.9
100	10M	0.7	98.6	98	2.0	414	6.1	120	2.9
1	100M	610.0	1.0	1342	1.0	stop	NA	2932	1.0
10	100M	56.2	10.9	210	6.4	2594	NA	432	6.8
100	100M	5.7	107.8	224	6.0	920	NA	220	13.3

times are closest to the theoretical maximum speedup. Intuitively, 1 thread cannot interleave the intensive CPU computation for  $\Gamma$  with I/O to scan  $X$  and at the other extreme 4 threads incur in too much overhead switching thread context and interleaving scan access to different chunks with one disk. Evidently, it is not worth going beyond the number of cores in the CPU, 4 in this case (results for 8 threads are much worse). In conclusion, the only alternative to achieve better speedup is to increase  $N$  in the parallel cluster (i.e., scale out).

We conclude our study of parallel processing analyzing speedup and scalability in a large computer cluster in the Amazon cloud. Since  $\Gamma$  is a summarization matrix whose computation is highly parallel it is natural to wonder if it works well on another

parallel system, different from SciDB. The current trend in big data analytics is Apache Spark [34], which has proven to be easier to program, more scalable and more general than MapReduce. With such motivation, we studied how to compute the same machine learning models on Spark. Spark has a rich library of linear algebra and machine learning functions called MLlib, which offers matrix multiplication, PCA and LR, but not VS. After trying several alternatives to compute  $\Gamma$  we picked the Gramian() function because it was the most efficient, which multiplies a matrix by itself, using  $Z$  as input. We should stress that according to Spark documentation Gramian() calls ScaLAPACK to evaluate the matrix multiplication in parallel in the most efficient manner. On the other hand, the PCA and LR models were directly available in MLlib passing the data set as the main input. The main difference between them is that PCA is solved by SVD like ours, but LR uses a gradient descent method [12, 16], which today is the fastest method to compute machine learning models. That is, we compare two very different approaches to analyze big data: our summarization matrix versus gradient descent.

Before discussing results we must explain installation, configuration and tuning of both systems. We had to install SciDB as a collection of SciDB instances, where we assigned 2 instances (threads) to each node. One of the nodes served both as coordinator and worker. That is, this node was “overworked”. Such limitation was an architecture constraint of SciDB since it is assumed that all nodes are uniform and the coordinator incurs on minimum overhead. On the other hand, Spark had the advantage of having a separate coordinator node, which had more RAM as shown in Table 6.2. Another important aspect was using sparse or dense matrices. Since Spark offered algorithms only with a dense (row) storage that is the one we used

in SciDB; time differences would be even larger processing the KDDnet data set in sparse form. Finally, we used the parameter defaults in Spark for LR and PCA models (20 iterations for LR, 5 components for PCA).

Table 6.8 compares SciDB and Spark running in the cloud varying  $N$  (the number of nodes) and  $n$  (data set size). The first major observation is that in R+SciDB the computation of  $\Theta$  (the model) taking  $\Gamma$  as input, took less than 1 second for LR and PCA. Therefore, we report the three times together in one column for SciDB. From a raw time perspective, SciDB was two orders of magnitude faster than Spark to compute  $\Gamma$ . To compute models the time difference becomes even more significant, close to three orders of magnitude, because the actual model computation took  $< 1$  sec. In summary, R+SciDB exhibited optimal speedup and significantly better performance. From a speedup perspective, to compute  $\Gamma$  Spark had good speedup at  $N = 10$ , but it was bad at  $N = 100$  (explained by a slow matrix transposition). To compute LR there was indeed better speedup up to  $N = 100$ , but far from ideal. Finally, for PCA results were the best, with decent speedup up to  $N = 100$ . On the other hand, SciDB exhibited ideal speedup to compute  $\Gamma$  both at  $N = 10$  and  $N = 100$ . Notice that due to SciDB architecture, when  $N = 1$  a single node is overworked (coordinator+worker), which explains why speedup is better than linear when  $N$  grows. Since the time to compute the LR and PCA models was negligible given the fast multicore CPU the speedup was the same for  $\Gamma$ , LR and PCA.

# Chapter 7

## Conclusions

We introduced Gamma, a comprehensive, but small, summarization matrix, which accelerates the computation of many machine learning models. We developed parallel summarization algorithms based on a special form of matrix multiplication to efficiently compute Gamma. A key optimization was evaluating matrix multiplication via a sum of vector outer products, instead of computing a product between the input matrix and its transpose. Based on the summarization matrix, algorithms to compute linear models are transformed to work in two phases: Summarization and iteration. Considering that dense and sparse matrices have different storage and require different processing in RAM, we introduced specialized summarization algorithms for dense and sparse matrices, respectively. We explained how to implement our algorithms in the SciDB array DBMS, to remove main memory limitations from R and enable parallel processing. We presented an extensive experimental section. R is significantly accelerated by Gamma and our algorithms working on SciDB eliminate R main memory limitations. Moreover, the array DBMS is an order of magnitude faster than a fast column-based DBMS, which evaluates matrix multiplication with SQL queries. We also compared with the fast matrix multiplication functions from LAPACK (MKL), the math library internally called by R. LAPACK (and MKL) fails when the input matrix does not fit in RAM, whereas our algorithms scale beyond main memory limits. Our algorithms are slower than LAPACK with small dense input matrices, but become significantly faster with large sparse matrices. Further benchmarking comparing the array DBMS with Spark on a large parallel cluster in



the cloud shows our algorithms running on the array DBMS are two orders of magnitude faster than Spark calling functions from MLlib. From a parallel perspective, Spark did not have good speedup to compute the summarization matrix (with a Gramian product via LAPACK), but had good (almost linear) speedup to compute PCA and LR. On the other hand, our algorithms running on the array DBMS showed linear speedup for both summarization and model computation.

Our experiments provide evidence that data summarization with fast matrix multiplication is a promising research direction. We need to study how other machine learning models, which require matrix multiplications, can exploit or generalize our Gamma summarization matrix. Large, high dimensional, data sets, are generally sparse matrices. Therefore, they deserve more attention than lower dimensional, dense, data sets. Based on current hardware trends, another reasonable assumption we made is that the partial and result matrices fit in main memory at each processing node. Nevertheless, an ideal summarization operator should work with input and output matrices of unlimited size. Clearly, our summarization matrix is promising to analyze big data, where we showed our array-based algorithms vastly outperform Spark, the dominating big data system. This big gap motivates integrating our summarization algorithms with Spark. We envision two alternatives: via a low-level transfer interface between the array DBMS and Spark, or programming them in Scala or Java.

# Bibliography

- [1] R. Agrawal, T. Imielinski, and A. Swami. Database mining: A performance perspective. *IEEE TKDE*, 5(6):914–925, 1993.
- [2] K. Akbudak and C. Aykanat. Simultaneous input and output matrix partitioning for outer-product-parallel sparse matrix-matrix multiplication. *SIAM J. Scientific Computing*, 36(5), 2014.
- [3] E. Anderson, Z. Bai, J. Dongarra, A. Greenbaum, A. McKenney, J. D. Croz, S. Hammerling, J. Demmel, C. Bischof, and D. Sorensen. Lapack: A portable linear algebra library for high-performance computers. In *Proc. of ACM/IEEE Supercomputing Conference*, pages 2–11, 1990.
- [4] P. Baumann, A. M. Dumitru, and V. Merticariu. The array database that is not a database: file based array query answering in rasdaman. In *Advances in Spatial and Temporal Databases*, pages 478–483. Springer, 2013.
- [5] P. Bradley, U. Fayyad, and C. Reina. Scaling clustering algorithms to large databases. In *Proc. ACM KDD Conference*, pages 9–15, 1998.
- [6] A. Buluç and J. Gilbert. Parallel sparse matrix-matrix multiplication and indexing: Implementation and experiments. *SIAM J. Scientific Computing*, 34(4), 2012.
- [7] J. Demmel. *Applied Numerical Linear Algebra*. SIAM, 1st edition, 1997.
- [8] J. Dongarra, I. Duff, D. Sorensen, and H. van der Vost. *Numerical Linear Algebra for High-Performance Computers*. SIAM, 1998.
- [9] P. Drineas, M. Mahoney, and S. Muthukrishnan. Sampling algorithms for l2 regression and applications. In *Proc. SODA*, pages 1127–1136, 2006.
- [10] W. DuMouchel, C. Volinski, T. Johnson, and D. Pregybon. Squashing flat files flatter. In *Proc. ACM KDD Conference*, 1999.
- [11] A. Gelman, J. Carlin, H. Stern, and D. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2003.
- [12] R. Gemulla, E. Nijkamp, P. Haas, and Y. Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In *Proc. KDD*, pages 69–77, 2011.
- [13] J. Han and M. Kamber. *Data Mining: Concepts and Techniques*. Morgan Kaufmann, San Francisco, 2nd edition, 2006.

- [14] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, New York, 1st edition, 2001.
- [15] J. Hellerstein, C. Re, F. Schoppmann, D. Wang, and et. al. The MADlib analytics library or MAD skills, the SQL. *Proc. of VLDB*, 5(12):1700–1711, 2012.
- [16] S. Hoi, J. Wang, P. Zhao, R. Jin, and P. Wu. Fast bounded online gradient descent algorithms for scalable kernel-based online learning. In *Proc. ICML*, 2012.
- [17] V. Karakasis, T. Gkountouvas, K. Kourtis, G. Goumas, and N. Koziris. An extended compression format for the optimization of sparse matrix-vector multiplication. *IEEE Trans. Parallel Distrib. Syst.*, 24(10):1930–1940, 2013.
- [18] F. Li and S. Nath. Scalable data summarization on big data. *Distributed and Parallel Databases*, 32(3):313–314, 2014.
- [19] J. Marin and C. Robert. *Bayesian Core: A Practical Approach to Computational Bayesian Statistics*. Springer, 2007.
- [20] A. Menon and C. Elkan. Fast algorithms for approximating the singular value decomposition. *ACM TKDD*, 5(2):13, 2011.
- [21] C. Ordonez. Integrating K-means clustering with a relational DBMS using SQL. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 18(2):188–201, 2006.
- [22] C. Ordonez. Optimization of linear recursive queries in SQL. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 22(2):264–277, 2010.
- [23] C. Ordonez. Statistical model computation with UDFs. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 22(12):1752–1765, 2010.
- [24] C. Ordonez, Y. Zhang, and W. Cabrera. The Gamma operator for big data summarization on an array DBMS. *Journal of Machine Learning Research (JMLR): Workshop and Conference Proceedings (BigMine 2014)*, 36:61–96, 2014.
- [25] C. Ordonez, Y. Zhang, and W. Cabrera. The Gamma matrix to summarize dense and sparse data sets for big data analytics. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2016.
- [26] M. Stonebraker, D. Abadi, A. Batkin, X. Chen, M. Cherniack, M. Ferreira, E. Lau, A. Lin, S. Madden, E. O’Neil, P. O’Neil, A. Rasin, N. Tran, and S. Zdonik. C-Store: A column-oriented DBMS. In *Proc. VLDB Conference*, pages 553–564, 2005.

- [27] M. Stonebraker, D. Abadi, D. DeWitt, S. Madden, E. Paulson, A. Pavlo, and A. Rasin. MapReduce and parallel DBMSs: friends or foes? *Commun. ACM*, 53(1):64–71, 2010.
- [28] M. Stonebraker, P. Brown, D. Zhang, and J. Becla. SciDB: A Database Management System for Applications with Complex Analytics. *Computing in Science and Engineering*, 15(3):54–62, 2013.
- [29] M. Stonebraker, S. Madden, D. J. Abadi, S. Harizopoulos, N. Hachem, and P. Helland. The end of an architectural era: (it’s time for a complete rewrite). In *VLDB*, pages 1150–1160, 2007.
- [30] M. Stonebraker, L. Rowe, and M. Hirohama. The implementation of Postgres. *IEEE TKDE*, 2(1):125–142, 1990.
- [31] R. Xin, J. Rosen, M. Zaharia, M. Franklin, S. Shenker, and I. Stoica. Shark: SQL and rich analytics at scale. In *Proc. ACM SIGMOD*, pages 13–24, 2013.
- [32] H. Xiong, S. Shekhar, P. Tan, and V. Kumar. TAPER: A two-step approach for all-strong-pairs correlation query in large databases. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 18(4):493–508, 2006.
- [33] A. Yzelman and D. Roose. High-level strategies for parallel shared-memory sparse matrix-vector multiplication. *IEEE Trans. Parallel Distrib. Syst.*, 25(1):116–125, 2014.
- [34] M. Zaharia, M. Chowdhury, M. Franklin, S. Shenker, and I. Stoica. Spark: Cluster computing with working sets. In *HotCloud USENIX Workshop*, 2010.
- [35] T. Zhang, R. Ramakrishnan, and M. Livny. BIRCH: An efficient data clustering method for very large databases. In *Proc. ACM SIGMOD Conference*, pages 103–114, 1996.
- [36] Y. Zhang, C. Ordonez, and W. Cabrera. Big data analytics integrating a parallel columnar DBMS and the R language. In *Proc. of IEEE CCGrid Conference*, 2016.