

Adaptive Consistency Management for Distributed Machine Learning

by

Christoph Alt

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Master of Science in Computer Science and Engineering

at the

TECHNISCHE UNIVERSITÄT BERLIN

April 2017

Author
Department of Electrical Engineering and Computer Science
April 15, 2017

Certified by.....
Prof. Dr. Odej Kao
Associate Professor
Thesis Supervisor

Certified by.....
Prof. Dr. Voker Markl
Associate Professor

Adaptive Consistency Management for Distributed Machine Learning

by

Christoph Alt

Submitted to the Department of Electrical Engineering and Computer Science
on April 15, 2017, in partial fulfillment of the
requirements for the degree of
Master of Science in Computer Science and Engineering

Abstract

In recent years, machine learning emerged as one of the most important parts of many successful applications and businesses. While the rise of artificial intelligence continues, the amount of data to be processed grows even faster. Unsurprisingly a lot of research has since focussed on methods to parallelize commonly used algorithms and new systems and frameworks have been released to improve the performance of distributed machine learning. Even though there has been a lot of progress towards more efficient systems, many of the state-of-the-art systems still have limitations.

Current frameworks are neither expressible nor flexible enough to allow efficient development of distributed machine learning algorithms, which makes them unsuited for experimentation and quick prototyping even though this is an essential part to optimize performance. On the other hand, most parallelization schemes exploit the algorithms stochastic nature to allow for parallel execution at the expense of lowered consistency among the distributed data structures. Even though this does not necessarily affect quality of the results, an improper level of consistency can severely affect algorithm performance, resulting in a non optimal convergence rate and therefore increased runtime.

The thesis introduces a novel framework for distributed machine learning, which is using a state centric programming model with yield semantics. This programming model allows the user to focus on the key parts of developing a distributed machine learning algorithm, namely update communication, update merging and consistency management among distributed workers while the system takes care of distributing the computation in an optimal fashion. The experiments show increase performance compared to a state-of-the-art data processing system like Apache Spark and at the same time reduce the effort of developing and experimenting with distributed machine learning algorithms at scale. What were the key findings or results? State of the art results, rapid prototyping and comparison of synchronization schemes

What is the significance or implications of the results?

Thesis Supervisor: Prof. Dr. Odej Kao
Title: Associate Professor

Acknowledgments

This is the acknowledgements section. You should replace this with your own acknowledgements.

Contents

1	Introduction	11
1.1	MapReduce and Beyond	12
1.2	Distributed Machine Learning	13
1.3	Thesis Outline	14
2	Background	17
2.1	Algorithms and Optimization	17
2.1.1	Iterative Convergent Algorithms	17
2.1.2	Optimization	18
2.1.3	CoCoA	20
2.2	Dataflow Systems	20
2.3	Distributed Machine Learning	20
2.3.1	Parameter Server	21
2.3.2	Consistency	21
3	State Centric Programming Model	23
4	Consistency Management	25
5	Experiments	27
6	Conclusion	29

List of Figures

2-1	Parameter Server	21
-----	----------------------------	----

Eidesstattliche Erklärung Hiermit erkläre ich Eides statt, dass ich dir vorliegende Arbeit selbstständig und ohne unerlaubte fremde Hilfe angefertigt, andere als die angegebenen Quellen und Hilfsmittel nicht benutzt und die den benutzten Quellen wörtlich oder inhaltlich entnommenen Stellen als solche kenntlich gemacht habe.

Unterschrift Ort, Datum

Chapter 1

Introduction

One of the most challenging tasks in computer science and engineering resolves around improving algorithm performance. In general this has been done by making hardware faster and inventing new strategies and algorithms to parallelize work more efficiently. Since it is clear that Moore's Law will not hold anymore, a lot of effort has been spent to horizontally scale algorithm computation across multiple machines. Machine learning is no exception and efficient parallelization is a key aspect towards more intelligent systems. By now, many general purpose frameworks for large scale data processing have been developed and published. Many of those are used for running more complex machine learning algorithms at scale as well. Unfortunately, the performance often is not satisfying due to the architecture and programming model not reflecting the underlying structure of most commonly used machine learning algorithms. Common data processing tasks can be represented as an extract-transform-load (ETL) pipeline, which is often easily parallelizable. This does not hold for machine learning, where algorithms are mostly sequential in nature and the only way of enabling parallel computation is by exploiting their inherent stochastic properties. This allows to break the sequential execution in favor of parallel learning on subsets, which then needs to be combined in order to obtain a global solution to the task. While this can lead to a great speedup in terms of the amount of data processed, it can have a negative affect on the learning progress. Therefore a key part of horizontally scaling machine learning algorithms is to ensure all participating learners have a consistent

view on each others progress while at the same time maintaining a tradeoff between communicating progress and spending time on their own learning task.

1.1 MapReduce and Beyond

Many of todays successful businesses throughout fields like finance, e-commerce and healthcare rely heavily on the ability to process vast amounts of user or sensorical data, collected to make services smarter and user experience better. In order to learn from the collected data, discover patterns and ultimately gain new insights, it needs to be processed by an algorithm. It is not uncommon that the input ranges somewhere between hundred gigabytes and tenth of petabytes. In the past, processing this much data required either a supercomputer, which was only available to large institutions or government entities or some proprietary compute cluster. All this changed when Google introduced MapReduce [1] in 2004. The MapReduce framework made it possible to process data in a distributed and fault tolerant way with the help of a compute cluster formed by hundredth or thousandth of machines. Instead of using a single, expensive, special hardware supercomputer the framework provides the foundation to assemble commodity hardware machines into a compute cluster. The framework takes care of all necessary aspects to ensure a fault tolerant and parallel execution of a task submitted to the cluster. The advantage compared to previous approaches is that the framework can be run entirely on top of machines using commodity hardware, which does not require special hardware and therefore equals low cost.

MapReduce essentially led the path to a convenient and widespread use of big data processing, which found its open source implementation in the Apache Hadoop project [2]. The project quickly gained traction and has spawned many business grade platforms, which quickly gained widespread adoption and by now provides a whole ecosystem around big data processing. The software stack includes a fault tolerant distributed filesystem (HDFS) a MapReduce framework and a cluster resource manager (YARN) [3]. On the other hand, MapReduce suffers from some practical

limitations that lead to the development of new, more sophisticated and specialised big data frameworks. With the most widely used frameworks being Apache Spark [4], Apache Flink [5] and GraphLab [6]. The first two frameworks use at its core a dataflow pipeline based architecture, whileas the latter uses a graph abstraction to model particular algorithms. All this works well for algorithms that can be expressed as an extract-transform-load (ETL) pipeline and are often embarrassingly parallel in nature. On the other hand, machine learning algorithms often rely heavily on many, computationally light, iterations to iteratively update a shared state (model) such as logistic regression or latent dirichlet allocation (LDA). These so called iterative-convergent algorithms required a change in how systems for distributed machine learning operate at its core.

1.2 Distributed Machine Learning

This limitation essentially sparked the development of specialized frameworks for distributed execution of iterative-convergent algorithms used in common machine learning tasks. The most widely recognized systems are Petuum [7], ParameterServer [8] and MALT [9]. Different from the MapReduce paradigm, these frameworks, instead of using a pipeline to transform an immutable dataset into another immutable dataset, operate on a fixed but mutable state, which is held by a single machine or distributed over multiple machines. This state can then be updated by workers that have computed an update locally and send it to these state keepers. Which act as surrogates, taking state updates and incorporating these into the state by some user defined function. While these systems can increase the performance on machine learning algorithms by an order of magnitude [7] compared to dataflow systems, most systems come with either limited usability, which makes it difficult to implement additional algorithms, are tied to a specific algorithm or are very low level frameworks. Efficiently distributing machine learning algorithms while at the same time provide the ability to conscisely express machine learning algorithms remains an extremely challenging problem. A system targeting the execution of those algorithms at scale must

therefore provide the ability to concisely express the underlying algebraic structure and at the same time be flexible enough to allow experimentation and fine tuning. Where consistency management is an essential part to ensure that algorithms are executed both fast and also learning performance is well enough.

1.3 Thesis Outline

In this thesis we introduce a novel framework for large scale distributed machine learning. It improves upon currently available systems by providing a powerful programming abstraction that can concisely express state of the art machine learning algorithms and at the same time minimizes the effort necessary to move from a single machine to a cluster. The framework design is optimized for efficient parallel asynchronous execution of iterative-convergent algorithms in a cluster and ensures the required consistency is enforced among parallel learners, depending on the algorithm properties. By allowing the user to decide how to maintaining the best tradeoff between algorithm execution and progress communication the performance is improved as well. All of this can be easily customized for quick prototyping and finetuning, which makes the system suited for developers as well as researchers. The goal of this thesis is to implement the state centric programming model and show its performance in comparison with Apache Spark on an example implementation of the CoCoA [10] framework. I start off by providing a background on the architecture and inner workings of state of the art frameworks for big data processing and distributed machine learning in Chapter 2. Furthermore the most commonly used algorithms and optimization techniques are introduced. The majority of those algorithms can be classified into the group of so called iterative-convergent algorithms for which I am going to provide a more formal description. I introduce the common algorithm parallelization strategies in distributed machine learning. I conclude the chapter by providing an overview over the challenges and issues that need to be addressed and considered when developing a distributed machine learning system and how this is achieved by current frameworks. This will give rise to understanding why a differ-

ent framework architecture and abstraction is necessary to improve the performance and expressability of large scale distributed machine learning applications. Chapter 3 therefore introduces the state centric programming model, which treats the state as a mutable first class citizen, which can be distributed and altered by updates that result from distributed computation. This allows the system to reason about the most optimal distribution of state in the cluster which is then scheduled with computation that can update the state. I further describe the architecture of the system and how the essential components are implemented. When updating a shared state from multiple different locations, consistency must be maintained in order to ensure the algorithm progresses as expected. The system is responsible for managing a state's consistency among all of its instances across the cluster. Chapter 4 therefore describes several schemes for ensuring consistency and at the same time optimizing for bandwidth and computational resource usage. In order to show the system and its consistency management in action, Chapter 5 compares the system against Apache Spark by running the CoCoA framework on two datasets with linear regression as the chosen algorithm. Chapter 6 summarizes the experiments with the lessons learned and provides suggestions on how to further improve systems for large scale distributed machine learning.

Chapter 2

Background

The section should provide the reader with enough background to follow the arguments in the following chapters regarding state centric programming model and consistency management. This includes an understanding of algorithms and optimization techniques commonly used in practice (not limited to distributed machine learning), the current state-of-the-art in dataflow systems and how dataflow is used to provide a fault tolerant and distributed framework for large scale data processing and machine learning. Furthermore the field of distributed machine learning is explained in more detail, including the current state-of-the-art frameworks used for this purpose, their limitations and what challenges arise when machine learning algorithms are executed in a distributed fashion on multiple machines.

2.1 Algorithms and Optimization

2.1.1 Iterative Convergent Algorithms

Consider a supervised learning setup with a dataset $D = \{z_1, \dots, z_n\}$ with each example z_i being represented by a pair (x_i, y_i) consisting of an input x_i and a scalar output y_i . Consider also a loss function $\ell(\hat{y}, y)$ quantifying the cost of predicting \hat{y} when the true output is y . As a model, a family F of functions $f_w(x)$ parameterized by a weight vector w is chosen. The goal is to find a function $f \in F$ that minimizes the

loss $Q(z, w) = \ell(f_w(x), y)$. Empirical risk $E_n(f) = \frac{1}{n} \sum_{i=0}^n \ell(f(x_i), y_i)$ performance on training set, expected risk generalization performance.

$$E_n(f_w) = \frac{1}{n} \sum_{i=0}^n \ell(f_w(x_i), y_i) \quad (2.1)$$

In order to find an optimal solution many algorithms used in large scale machine learning such as regression, topic models, matrix factorization or neural networks employ either gradient based methods or markov chain monte carlo methods. To obtain the optimal solution those algorithms try to iteratively update the weight vector w . At each iteration t an updated weight vector w^t is computed based on the vector of the previous iteration $w^{(t-1)}$ and the data D . The resulting model f_{w^t} is again a better summary of the data D under the objective Q . Eq. 2.2 shows the process of refining the model, with Δ being an arbitrary update function.

$$w^t = w^{(t-1)} + \Delta(w^{(t-1)}, D) \quad (2.2)$$

The update function depends on the algorithm employed and can be viewed as a procedure of obtaining a step towards a better model. At each iteration an update Δw is computed and applied to the previous weight vector until a stopping condition is satisfied. E.g. the distance to the optimal solution or the objective difference between two iterations is monitored. When the difference is below a certain threshold the computation stops and the algorithm is said to be converged.

2.1.2 Optimization

In order to estimate the optimal parameters w^* of a function belonging to class $f_{w^*} \in F$, numerous techniques can be employed to estimate said parameters. In many cases, especially large scale machine learning methods such as (stochastic) gradient descent and coordinate ascent are used to iteratively optimized the parameterization of the chosen function class. Both techniques represent different rules of computing the update shown in 2.2. Gradient descent updates the weights w at each iteration t

on the basis of the gradient of $E_n(f_w)$,

$$w^t = w^{(t-1)} - \eta \frac{1}{n} \sum_{i=0}^n \nabla_w Q(z_i, w^{(t-1)}) \quad (2.3)$$

where η is a chosen gain, often referred to as learning rate. While this achieves linear convergence under sufficient regularity assumptions and a sufficiently small learning rate η [11] [12] a more simplified version is commonly used in practice, called stochastic gradient descent (SGD). Instead of computing the gradient $\nabla_w E_n(f_w)$ exactly, the gradient is estimated at each iteration t based on a single randomly picked example z_t .

$$w^t = w^{(t-1)} - \eta_t \nabla_w Q(z_t, w^{(t-1)}) \quad (2.4)$$

The assumption is that the gradient obtained by 2.4 behaves similar to its expectation in 2.3. The convergence properties have been studied extensively and under mild conditions an almost sure convergence can be established when the learning rate satisfies the conditions $\sum_t \eta_t^2 < \infty$ and $\sum_t \eta_t = \infty$ [12]. The general structure of stochastic gradient descent is described in Algorithm 1.

Algorithm 1 Stochastic Gradient Descent

- 1: $k \leftarrow 1$ and initialize $w^0 \in \mathbb{R}^d$
 - 2: **repeat**
 - 3: **for do**
 - 4: $w^t \leftarrow w^{(t-1)} - \eta_t \nabla_w Q(z_t, w^{(t-1)})$
 - 5: **until** termination criteria satisfied
-

Coordinate descent on the other hand iteratively tries to optimize a given objective by successively performing approximate minimization along a coordinate direction while keeping the other directions fixed.

Algorithm 2 Stochastic Coordinate Ascent

- 1: $k \leftarrow 1$ and initialize $w^0 \in \mathbb{R}^d$
 - 2: **repeat**
 - 3: **for do**
 - 4:
 - 5: **until** termination criteria satisfied
-

2.1.3 CoCoA

Due to their widespread application in large scale machine learning and recent advances in the field of distributed optimization the thesis focuses on linear regularized objectives. One of the advances is the publication of a framework for distributed optimization called CoCoA (Communication-efficient distributed dual Coordinate Ascent) [10], which leverages the primal-dual structure of the beforementioned optimization problems. The proposed strategy helps effectively combining the results from local computation without having to deal with conflicts resulting from similar updates computed on other machines. The framework supports objectives of the form

$$Q(z, w) = \ell(f_w(x), y) + r(w) \quad (2.5)$$

where ℓ is convex and smooth and r is assumed to be separable. In this context separable means $r(x) = \sum_{i=0}^n r_i(x_i)$. Commonly the term ℓ is an empirical loss over the data of the form $\sum_i \ell(f_w(x_i), y_i)$ and the term r is a regularizer, e.g. $r(w) = \lambda \|w\|_p$ where λ is a regularization parameter. Many of the beforementioned algorithms in machine learning can be expressed in this form, such as logistic and linear regression, lasso and sparse logistic regression and support vector machines.

2.2 Dataflow Systems

2.3 Distributed Machine Learning

Distributed machine learning poses a number of unique challenges resulting from the fact that most iterative convergent algorithms are sequential in nature. Parallelizing such an algorithm can be done by exploiting either the stochastic nature inherent to these algorithms or by exploiting decomposability of model parameters. Doing so leaves the developer with numerous possibilities and responsibilities regarding the distribution of data across the workers, the formation of workers and the communication and synchronization scheme to be employed. Various frameworks have been

published which all tackle these problems in a specific way. Often focusing on a fixed master/worker architecture with specific algorithm.

List the challenges of distributed ML (synchronization/consistency, distributing state and model, choosing the best formation of workers (maybe different jobs (computing updates, merging updates))), choosing the best communication model, show picture

2.3.1 Parameter Server

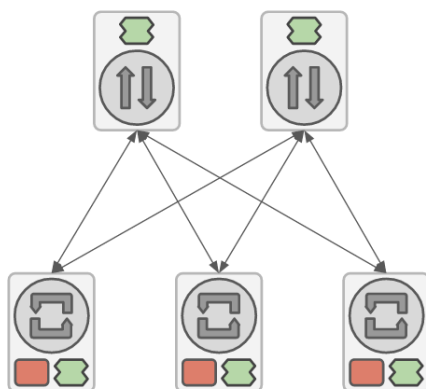


Figure 2-1: Parameter Server

2.3.2 Consistency

The most important part of any distributed system is the synchronization strategy used to ensure consistency among multiple nodes concurrently accessing and updating the parameters stored on the parameter server. There are multiple schemes to synchronize nodes during the iterative parameter refinement. *Bulk synchronous parallelization (BSP)* leads to the best algorithm throughput (e.g. convergence achieved over the number of data points processed). Essentially each worker must finish its iteration and push all updates to the parameter server. The server then computes a refined model according to Eq. 2.6 and each node retrieves the updated parameters before beginning the next iteration. This synchronization scheme guarantees

consistency among all nodes at all times.

$$W^t = W^{(t-1)} + \frac{1}{K} \sum_{k=1}^K \Delta(W_k^{(t-1)}, D_k) \quad (2.6)$$

While this synchronization strategy essentially recovers the sequential algorithm for a single machine and has the same convergence properties and guarantees, it suffers from a severe limitation when used in a distributed setup [?]. Imagine one of the workers is for some reason a lot slower than the others. Due to the synchronization strategy, the other workers have to wait for this particular worker to complete its iteration. This is well known as the straggler problem [?] and can seriously affect performance in a distributed environment, because the progress is limited by the slowest node in the cluster. The second strategy is *total asynchronous parallelization (TAP)*. Similar to BSP, all nodes push their parameter updates to the server after each iteration but in this case, the changes are applied to the model immediately. No waiting for other workers is required, resulting in a very high data throughput. The straggler problem can be mitigated by this synchronization scheme as well, depicted in Figure ?? . Even though worker 3 is a straggler, the remaining workers can continue with their next iterations without waiting for slower workers. Although this consistency scheme seems to work quite well in practice [?], it lacks formal convergence guarantees and can even diverge [?]. The reason is that no bound exists for a situation where the divergence in iterations between the slowest and the fastest worker is unbound. A middle ground between bulk synchronous parallelization and total asynchronous parallelization is *stale synchronous parallel (SSP)* [?] or *bounded staleness (BS)*. As shown in Figure ?? , BSP introduces a maximum delay, or staleness threshold, of Δ_{max} between the slowest and fastest node. This overcomes the limitation of the TAP approach by introducing a bound on the number of iterations. Formal convergence guarantees can be restored while still maintaining the flexibility of asynchronous parallelization and limiting the straggler problem [?].

Chapter 3

State Centric Programming Model

Chapter 4

Consistency Management

Chapter 5

Experiments

Chapter 6

Conclusion

Bibliography

- [1] J. Dean and S. Ghemawat, “MapReduce: Simplified Data Processing on Large Clusters,” *Proceedings of 6th Symposium on Operating Systems Design and Implementation*, pp. 137–149, 2004.
- [2] A. Hadoop, “Hadoop,” 2009.
- [3] V. Kumar Vavilapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth, B. Saha, C. Curino, O. O ’malley, S. Radia, B. Reed, and E. Baldeschwieler, “Apache Hadoop YARN: Yet Another Resource Negotiator,” *SOCC ’13 Proceedings of the 4th annual Symposium on Cloud Computing*, vol. 13, pp. 1–3, 2013.
- [4] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, “Spark : Cluster Computing with Working Sets,” *HotCloud’10 Proceedings of the 2nd USENIX conference on Hot topics in cloud computing*, p. 10, 2010.
- [5] A. Alexandrov, R. Bergmann, S. Ewen, J. C. Freytag, F. Hueske, A. Heise, O. Kao, M. Leich, U. Leser, V. Markl, F. Naumann, M. Peters, A. Rheinländer, M. J. Sax, S. Schelter, M. Höger, K. Tzoumas, and D. Warneke, “The Stratosphere platform for big data analytics,” *VLDB Journal*, vol. 23, no. 6, pp. 939–964, 2014.
- [6] Y. Low, J. Gonzalez, A. Kyrola, D. Bickson, and C. Guestrin, “GraphLab: A Distributed Framework for Machine Learning in the Cloud,” *The 38th International Conference on Very Large Data Bases*, vol. 5, no. 8, pp. 716–727, 2012.
- [7] E. P. Xing, Q. Ho, W. Dai, J. K. Kim, J. Wei, S. Lee, and X. Zheng, “Petuum : A New Platform for Distributed Machine Learning on Big Data,” 2015.
- [8] M. Li, D. G. Andersen, J. W. Park, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, B.-y. Su, M. Li, D. G. Andersen, J. W. Park, A. J. Smola, and A. Ahmed, “Scaling Distributed Machine Learning with the Parameter Server,” *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, 2014.
- [9] H. Li, A. Kadav, E. Kruus, and C. Ungureanu, “MALT : Distributed Data-Parallelism for Existing ML Applications,” 2015.

- [10] M. Jaggi, V. Smith, M. Tak??, J. Terhorst, S. Krishnan, T. Hofmann, and M. Jordan, “Communication-efficient distributed dual coordinate ascent,” *Advances in Neural Information Processing Systems*, vol. 4, no. January, 2014.
- [11] J. E. Dennis Jr and R. B. Schnabel, *Numerical methods for unconstrained optimization and nonlinear equations*. SIAM, 1996.
- [12] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proceedings of COMPSTAT’2010*, pp. 177–186, 2010.