

Templates for scalable data analysis

2 Synchronous Templates

Amr Ahmed, Alexander J Smola, Markus Weimer

Yahoo! Research & UC Berkeley & ANU















































































MAGIC Etch A Sketch® SCREEN Example Formation in Pig • Modeling today Hadoop, Spark, Pregel • Future Declarative Systems lonzertel Gid GING MARY "GENTIONE OF TORES The to be been used as a state of the second s



EMail

Log Files







































Requirements



Requirements









Apache Pig

- Relational Query Language
- Similar to SQL
- Performs runtime optimizations
- Executes Queries on Apache Hadoop
- Developed and heavily used by Yahoo!
- Open Source (Apache)



Pig: Example Formation

- Feature and Label Extraction
 - User Defined Function
 - Applied via FOREACH ... GENERATE
- Example formation
 - JOIN between the outputs of the above


- Parallel, Distributed programming framework
- User defines two functions:
 - map(x) emits (key, value) pairs
 - reduce(k, x[]) gets all values for a key, produces output















GroupBy (Shuffle)



Reduce























- Open Source MapReduce Implementation:
 - HDFS: Distributed FileSystem
 - YARN: Resource Management
 - MapReduce: Programming Framework

http://hadoop.apache.org



- Open Source MapReduce Implem Hadoop .23
 - HDFS: Distributed FileSystem
 - YARN: Resource Management
 - MapReduce: Programming Framework

http://hadoop.apache.org

New in

MapReduce for ML

- Learning algorithm can access the learning problem only through a statistical query oracle
- The statistical query oracle returns an estimate of the expectation of a function f(x,y) (averaged over the data distribution).

Efficient Noise-Tolerant Learning from Statistical Queries

MICHAEL KEARNS

AT&T Laboratories—Research, Florham Park, New Jersey

Abstract. In this paper, we study the problem of learning in the presence of classification probabilistic learning model of Valiant and its variants. In order to identify the class learning algorithms in the most general way, we formalize a new but related model of le *statistical queries*. Intuitively, in this model, a learning algorithm is forbidden to examine examples of the unknown target function, but is given access to an oracle providing probabilities over the sample space of random examples.

One of our main results shows that any class of functions learnable from statistical que learnable with classification noise in Valiant's model, with a noise rate approaching the theoretic barrier of 1/2. We then demonstrate the generality of the statistical query mo that practically every class learnable in Valiant's model and its variants can also be learned model (and thus can be learned in the presence of noise). A notable exception to this stat class of parity functions, which we prove is not learnable from statistical queries, and noise-tolerant algorithm is known.

Categories and Subject Descriptors: F. [Theory of Computation]; G.3 [Probability and St [Artificial Intelligence]; I.5 [Pattern Recognition]

General Terms: Computational learning theory, Machine learning

Additional Key Words and Phases: Computational learning theory, machine learning

1. Introduction

In this paper, we study the extension of Valiant's learning model [Val in which the positive or negative classification label provided with eac example may be corrupted by random noise. This extension was first ex the learning theory literature by Angluin and Laird [1988], who form simplest type of white label noise and then sought algorithms tole highest possible rate of noise. In addition to being the subject of a r theoretical studies [Angluin and Laird 1988: Laird 1988: Sloan 1988: K

MapReduce for ML

- Rephrase oracle in summation form.
- Map: Calculate function estimates over sub-groups of data.
- **Reduce:** Aggregate the function estimates from various sub-groups.

Map-Reduce for Machine Learning on Multicore

Cheng-Tao Chu * chengtao@stanford.edu

Sang Kyun Kim * skkim38@stanford.edu

Yi-An Lin * ianl@stanford.edu

YuanYuan Yu* yuanyuan@stanford.edu Gary Bradski^{*†} And garybradski@gmail ang@cs.

Andrew Y. Ng * ang@cs.stanford.edu

Kunle Olukotun * kunle@cs.stanford.edu

* CS. Department, Stanford University 353 Serra Mall, Stanford University, Stanford CA 94305-9025. † Rexee Inc.

Abstract

We are at the beginning of the multicore era. Computers will have increasingly many cores (processors), but there is still no good programming framework for these architectures, and thus no simple and unified way for machine learning to take advantage of the potential speed up. In this paper, we develop a broadly applicable parallel programming method, one that is easily applied to *many* different learning algorithms. Our work is in distinct contrast to the tradition in machine learning of designing (often ingenious) ways to speed up a *single* algorithm at a time. Specifically, we show that algorithms that fit the Statistical Query model [15] can be written in a certain "summation form," which allows them to be easily parallelized on multicore computers. We adapt Google's map-reduce [7] paradigm to demonstrate this parallel speed up technique on a variety of learning algorithms including locally weighted linear regression (LWLR), k-means, logistic regression (LR), naive Bayes (NB), SVM, ICA, PCA, gaussian discriminant analysis (GDA), EM, and backpropagation (NN). Our experimental results show basically linear speedup with an increasing number of processors.

1 Introduction

Frequency scaling on silicon—the ability to drive chips at ever higher clock rates—is beginning to hit a power limit as device geometries shrink due to leakage, and simply because CMOS consumes power every time it changes state [9, 10]. Yet Moore's law [20], the density of circuits doubling every generation, is projected to last between 10 and 20 more years for silicon based circuits [10].













- Machine Learning Library
- Implementations of many algorithms, both on Hadoop MapReduce and stand-alone
- Open Source (Apache)
- Welcoming, helpful community

http://mahout.apache.org



- Recommender Systems, e.g.
 - User and Item based recommenders
 - Collaborative Filtering
- Clustering (K-Means, Mean Shift, ...)
- Topic Models (LDA)
- Supervised ML
 - (Logistic) Regression
 - Linear SVMs
 - Decision Trees and Forests

Efficient Noise-Tolerant Learning from Statistical Oueries

MICHAEL KEARNS

Map-Reduce for Machine Learning on Multicore

AT&T Laboratories-Research, Florham Park, New Jersey

Yi-An Lin *

Abstract. In this paper, we study the problem of learning in the presence of classification noise in the probabilistic learning model of Valiant and its variants. In order to identify the class of "robust" learning algorithms in the most general way, we formalize a new but related model of learning from statistical queries. Intuitively, in this model, a learning algorithm is forbidden to examine individual examples of the unknown target function, but is given access to an oracle providing estimates of probabilities over the sample space of random examples.

One of our main results shows that any class of functions learnable from statistical queries is hr fact learnable with classification noise in Valiant's model, with a noise rate approaching the information theoretic barrier of 1/2. We then demonstrate the generality of the statistical outery model, showing that practically every class learnable in Valiant's model and its variant's can also be learned in the new model (and thus can be learned in the presence of noise). A notable exception to this statement is the class of parity functions, which we prove is not learnable from statistical queries, and for which no noise-tolerant algorithm is known.

Categories and Subject Descriptors: F. [Theory of Computation]; G.3 [Probability and Statistics]; I.2. [Artificial Intelligence]; I.5 [Pattern Recognition]

General Terms: Computational learning theory, Machine learning

Additional Key Words and Phases: Computational learning theory, machine learning

1. Introduction

In this paper, we study the extension of Valian in which the positive or negative classification example may be corrupted by random noise. This extension was first e the learning theory literature by Angluin and Laird [19 sistephism typeodelinghitten label moise and then sought algorithms tolerating the highest possible rate of noise. In addition to being the subject of a number of theoretical studies [Angluin and Laird 1988; Laird 1988; Sloan 1988; Kearns and Li 1993], the classification noise model has become a common paradigm for experimental machine learning search.

1 / 140 53

Tutorial @ KDD 2011 http://www.slideshare.net/hadoop

н () н



III second

Efficient Noise-Tolerant Learning from Statistical Oueries

MICHAEL KEARNS

Map-Reduce for Machine Learning on Multicore

AT&T Laboratories-Research, Florham Park, New Jersey

Yi-An Lin *

Abstract. In this paper, we study the problem of learning in the presence of classification noise in the probabilistic learning model of Valiant and its variants. In order to identify the class of "robust" learning algorithms in the most general way, we formalize a new but related model of learning from statistical queries. Intuitively, in this model, a learning algorithm is forbidden to examine individual examples of the unknown target function, but is given access to an oracle providing estimates of probabilities over the sample space of random examples.

One of our main results shows that any class of functions learnable from statistical queries is hr fact learnable with classification noise in Valiant's model, with a noise rate approaching the information theoretic barrier of 1/2. We then demonstrate the generality of the statistical outery model, showing that practically every class learnable in Valiant's model and its variant's can also be learned in the new model (and thus can be learned in the presence of noise). A notable exception to this statement is the class of parity functions, which we prove is not learnable from statistical queries, and for which no noise-tolerant algorithm is known.

Categories and Subject Descriptors: F. [Theory of Computation]; G.3 [Probability and Statistics]; I.2. [Artificial Intelligence]; I.5 [Pattern Recognition]

General Terms: Computational learning theory, Machine learning

Additional Key Words and Phases: Computational learning theory, machine learning

1. Introduction

In this paper, we study the extension of Valian in which the positive or negative classification example may be corrupted by random noise. This extension was first e the learning theory literature by Angluin and Laird [19 sistephism typeodelinghitten label moise and then sought algorithms tolerating the highest possible rate of noise. In addition to being the subject of a number of theoretical studies [Angluin and Laird 1988; Laird 1988; Sloan 1988; Kearns and Li 1993], the classification noise model has become a common paradigm for experimental machine learning search.

1 / 140 53

Tutorial @ KDD 2011 http://www.slideshare.net/hadoop

н () н



III second

Map-Reduce for Machine Learning on Multicore

Efficient Noise-Tolerant Learning from Statistical Queries		
MICHAEL CHARNS ATAL Laboration and Market No. Control *	Sang Kyun Kim *	Yi-An Lin *
chengtao@stanford.edu	skkim38@stanford.edu	ianl@stanford.edu
probabilistic learning model of Valiant and in variants. In order to identify the class of value learning algorithms in the most general way, se formalistic and the second of learning the second of the second sec	Gary Bradski * garybradski@gmail	Andrew Y.Ng* ang@cs.stanford.edu
General Terms: Computational learning theory, Machine learning Additional Key Words and Phases: Computational learning theory, machine learning	Kunle Olukotun *	
1. Introduction kur	le@cs.stanford.edu	
In this paper, we study the extension of Valiant's learning model (Valiant '9) in which the positive or negative classification label provided with each rande example may be corrupted by random nos. Perturn we want to reason simplest upper of which label noise and then sough argonithes to there are a study of the study of the sough argonithes to the study study (14) and and Laid 1985; Laint Ott Laint Li 1993], the classification noise model has become a common paradigm experimental machine learning research.	 stanford University 353 Se miversity, Stanford CA 94305-9 frequencies the Resee Inc. 	rra Mall, 0025.

Abstract



Map-Reduce for Machine Learning on Multicore

Efficient Noise-Tolerant Learning from Statistical Queries		
MICHAEL CHARNS ATAL Laboration and Market No. Control *	Sang Kyun Kim *	Yi-An Lin *
chengtao@stanford.edu	skkim38@stanford.edu	ianl@stanford.edu
probabilistic learning model of Valiant and in variants. In order to identify the class of value learning algorithms in the most general way, se formalistic and the second of learning the second of the second sec	Gary Bradski * garybradski@gmail	Andrew Y.Ng* ang@cs.stanford.edu
General Terms: Computational learning theory, Machine learning Additional Key Words and Phases: Computational learning theory, machine learning	Kunle Olukotun *	
1. Introduction kur	le@cs.stanford.edu	
In this paper, we study the extension of Valiant's learning model (Valiant '9) in which the positive or negative classification label provided with each rande example may be corrupted by random nos. Perturn we want to reason simplest upper of which label noise and then sough argonithes to there are a study of the study of the sough argonithes to the study study (14) and and Laid 1985; Laint Ott Laint Li 1993], the classification noise model has become a common paradigm experimental machine learning research.	 stanford University 353 Se miversity, Stanford CA 94305-9 frequencies the Resee Inc. 	rra Mall, 0025.

Abstract











Efficient Noise-Tolerant Learning from Statistical Queries

MICHAEL KEARNS AT&T Laboratories—Research, Florham Park, New Jersey

Abstract, In this paper, we study the problem of learning in the presence of closult-trains noise in the example algorithm in the most general way, we formulate a new but retained model of learning frame matrixed queries. The attractive, in this model, a learning algorithm is forbidden to constrained induced retained queries. The attractive, in this model, a learning algorithm is forbidden to constrained induced of the attractive, in this model, a learning algorithm is forbidden to constrained induced of the attractive, in this model, a learning algorithm the frame that the attractive in the learning of the attractive, in this model, and on frame time interaction frame training approximation of the attractive in the strategiest of the attractive in the attractive in the attractive in the attractive in the strategiest of the attractive in the strategiest of the attractive of the strategiest of the strategi

General Terms: Computational learning theory, Machine learning Additional Key Words and Phases: Computational learning theory, machine learning

1. Introduction

In this paper, we study the extension of Valiant's learning model [Valiant 1984], in which the positive or negative classification label provided with each random example may be corrupted by random noise. This extension was first examined in he learning theory fluentus: by Anguin and Laird [Valiss], who formalized the simplest type of white label noise and then sought algorithms tolerating the however the simplest of the simplest result of the simplest type of theoretical studies (Anguin and Laird 1988), Laird 1988. Shoan 1988, Scans 1988, Scans Li 1993), the classification noise model has become a common paradigm for exercimental matching learning research.

Map-Reduce for Machine Learning on Multicore

Cheng-Tao Chu * Sang Kyun Kim * Yi-An Lin * chengtao@stanford.edu skkim38@stanford.edu ianl@stanford.edu

YuanYuan Yu * Gary Bradski ⁴ Andrew Y.Ng * yuanyuan@stanford.edu garybradski@gmail ang@cs.stanford.edu

Kunle Olukotun * kunle∉cs.stanford.edu *CS. Department, Stanford University 353 Serra Mall, Stanford University, Stanford CA 94305-9025. ! Revee Inc.

Abstract

We are at the beginning of the multicose era. Computers will have increasingly many cores (processors), but there is still as pood programming framework for these architectures, and thus no simple and unified way for machine learning to plicable parallel programming methods (on the his cashy) applied to nany different learning of designing (often argosino) ways to speed or a particel algorithm at a time. Specifically, we show that algorithms that the Statistical Query model [15] and be written as accitant "manner body on synt so good on the statistical and party model at the synthemic accitant" statistical form," which allows them to be easily part and be written as accitant "manner body on the statistical Query model [15] at he written is accitant" statuse theory of the statistical Query model [16] and be written in accitant "manner body on the statistical Query model [16] at he written is accitant" statuse theory of the statistical Query model [16] at hey written in accitant "manner body on the statisticant" (16). A statistical Query model [16] at hey written in accitant" statuse theory of the statistical Query model [16] and the statisticant of the statisticant of the statisticant of the statisticant of the statisticant including locality weighted linear represens (ILWIR), hey mean mainter analysis (ICDA), EM, and backgropagation (NN), Our experimental results also basically



Tutorial @ KDD 2011 http://www.slideshare.net/hadoop



Trouble

- ML is iterative
- Each iteration is a Job
- Overhead per job (>45s)
 - Scheduling
 - Program Distribution
 - Data Loading and Parsing
 - State Transfer

Map-Reduce for Machine Learning on Multicore

Cheng-Tao Chu * chengtao@stanford.edu

Sang Kyun Kim * skkim38@stanford.edu

Yi-An Lin * ianl@stanford.edu

YuanYuan Yu * yuanyuan@stanford.edu Gary Bradski *†

Andrew Y.Ng * ang@cs.stanford.edu

Kunle Olukotun * kunle@cs.stanford.edu

garybradski@gmail

* CS. Department, Stanford University 353 Serra Mall, Stanford University, Stanford CA 94305-9025. † Rexee Inc.

Abstract

We are at the beginning of the multicore era. Computers will have increasingly many cores (processors), but there is still no good programming framework for these architectures, and thus no simple and unified way for machine learning to take advantage of the potential speed up. In this paper, we develop a broadly applicable parallel programming method, one that is easily applied to *many* different learning algorithms. Our work is in distinct contrast to the tradition in machine learning of designing (often ingenious) ways to speed up a *single* algorithm at a time. Specifically, we show that algorithms that fit the Statistical Query model [15] can be written in a certain "summation form," which allows them to be easily parallelized on multicore computers. We adapt Google's map-reduce [7] paradigm to demonstrate this parallel speed up technique on a variety of learning algorithms including locally weighted linear regression (LWLR), k-means, logistic regression (LR), naive Bayes (NB), SVM, ICA, PCA, gaussian discriminant analysis (GDA), EM, and backpropagation (NN). Our experimental results show basically linear speedup with an increasing number of processors.

1 Introduction

Frequency scaling on silicon—the ability to drive chips at ever higher clock rates—is beginning to hit a power limit as device geometries shrink due to leakage, and simply because CMOS consumes power every time it changes state [9, 10]. Yet Moore's law [20], the density of circuits doubling every generation, is projected to last between 10 and 20 more years for silicon based circuits [10].

Trouble

- ML is iterative
- Each iteration is a Job
- Overhead per job (>45s)
 - Scheduling
 - Program Distribution
 - Data Loading and Parsing
 - State Transfer



* CS. Department, Stanford University 353 Serra Mall, Stanford University, Stanford CA 94305-9025. ! Rexee Inc.

Abstract

We are at the beginning of the multicore era. Computers will have increasingly many cores (processors), but there is still no good programming framework for these architectures, and thus no simple and unified way for machine learning to take advantage of the potential speed up. In this paper, we develop a broadly applicable parallel programming method, one that is easily applied to *many* different learning algorithms. Our work is in distinct contrast to the tradition in machine learning of designing (often ingenious) ways to speed up a *single* algorithm at a time. Specifically, we show that algorithms that fit the Statistical Query model [15] can be written in a certain "summation form," which allows them to be easily parallelized on multicore computers. We adapt Google's map-reduce [7] paradigm to demonstrate this parallel speed up technique on a variety of learning algorithms including locally weighted linear regression (LWLR), k-means, logistic regression (LR), naive Bayes (NB), SVM, ICA, PCA, gaussian discriminant analysis (GDA), EM, and backpropagation (NN). Our experimental results show basically linear speedup with an increasing number of processors.

1 Introduction

Frequency scaling on silicon—the ability to drive chips at ever higher clock rates—is beginning to hit a power limit as device geometries shrink due to leakage, and simply because CMOS consumes power every time it changes state [9, 10]. Yet Moore's law [20], the density of circuits doubling every generation is projected to last between 10 and 20 more years for silicon based circuits [10].



Solutions

- Local (subsampling)
- MPI
- Spark
- Pregel
Subsampling

- Form examples on the cluster
- Subsample the data on the cluster
- Train a model on a single machine









Per-Partition Training

Averaging





Per-Partition Training

Averaging



Per-Partition Training



Message Passing Interface

- Mature HPC standard
- Supported on many clusters (e.g. OpenMPI)
- Available in C, Fortran and Scripting Languages
- Key operation here: AllReduce



... AllReduce()

... AllReduce()





... AllReduce()

... AllReduce()





... AllReduce()

... AllReduce()





... AllReduce()





... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()



... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()

... AllReduce()



Hadoop AllReduce

- Use Hadoop for
 - Data local scheduling
 - Good machine identification
- Use MPI for
 - AllReduce
- 30x Speedup over Hadoop MapReduce

A Reliable Effective Terascale Linear Learning System

Alekh Agarwal Department of EECS UC Berkeley alekh@cs.berkeley.edu

Miroslav Dudík Yahoo! Research New York, NY mdudik@yahoo-inc.com

ABSTRACT

We present a system and a set of techniques for learning linear predictors with convex losses on terascale datasets, with trillions of features,¹ billions of training examples and millions of parameters in an hour using a cluster of 1000 machines. Individually none of the component techniques is new, but the careful synthesis required to obtain an efficient implementation is a novel contribution. The result is, up to our knowledge, the most scalable and efficient linear learning system reported in the literature. We describe and thoroughly evaluate the components of the system, showing the importance of the various design choices.

1. INTRODUCTION

Distributed machine learning is a research area that has seen a growing body of literature in recent years. Much work focuses on problems of the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad \sum_{i=1}^n \ell(\mathbf{w}^\top \mathbf{x}_i; y_i) + \lambda R(\mathbf{w}),$$

where \mathbf{x}_i is the feature vector of the *i*-th example, y_i is the label, \mathbf{w} is the linear predictor, ℓ is a loss function and R a regularizer. Much of this work exploits the natural decomposability over examples in (1), partitioning the examples over different nodes in a distributed environment such as a cluster.

Perhaps the simplest learning strategy when the number of samples n is very large is to subsample a smaller set of examples that can be tractably learned with. However, this strategy only works if the problem is simple enough or the number of parameters is very small. The setting of interest here is when a large number of samples is really needed to learn a good model, and distributed algorithms are a natural choice for such scenarios.

¹The number of zero entries in th <u>http://hunch.net/~vw</u>

(1)

Olivier Chapelle Yahoo! Research Santa Clara, CA chap@yahoo-inc.com

John Langford Yahoo! Research New York, NY jl@yahoo-inc.com

Some prior works (McDonald et al., 2010; Zinkevich ϵ 2010) consider online learning with averaging and I et al. (2010a) propose gossip-style message passing rithms extending the existing literature on distributed vex optimization (Bertsekas and Tsitsiklis, 1989). Lan et al. (2009) analyze a delayed version of distributed o learning. Dekel et al. (2010) consider mini-batch versic online algorithms which are extended to delay-based up in Agarwal and Duchi (2011). A recent article of Boyd (2011) describes an application of the ADMM techniqu distributed learning problems. GraphLab (Low et al., ϵ is a parallel computation framework on graphs. More cl related to our work is that of Teo et al. (2007) who use I to parallelize a bundle method for optimization.

However, all of the aforementioned approaches see leave something to be desired empirically when deploye large clusters. In particular their throughput—measur the input size divided by the wall clock running tim smaller than the the I/O interface of a single machin almost all parallel learning algorithms (Bekkerman e 2011, Part III, page 8). The I/O interface is an upper b on the speed of the fastest sequential algorithm sinc sequential algorithms are limited by the network inte in acquiring data. In contrast, we were able to achie throughput of 500M features/s, which is about a factor faster than the 1Gb/s network interface of any one no

An additional benefit of our system is its compatiwith MapReduce clusters such as Hadoop (unlike MPI-b systems) and minimal additional programming effort tc allelize existing learning algorithms (unlike MapReduc proaches).

One of the key components in our system is a comr cation infrastructure that efficiently accumulates and be casts values across all nodes of a computation. It is func ally similar to MPI AllReduce (hence we use the name) it takes advantage of and is compatible with We doop so

MPI: Conclusion

• The Good

- Computational Performance
- Well established software available
- The Bad
 - No fault tolerance
- The Ugly
 - Ignorance of shared clusters
 - Systems-Level decisions at the algorithm layer

Spark: Intro

- Open Source cluster computation framework
- Developed at UC Berkeley by the AMP Lab
- Aimed at interactive and iterative use cases
- 30x faster than Hadoop for those
- User interface: Embedded Domain Specific
 Language in Scala

http://spark-project.org/

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
   val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
   w -= gradient
}</pre>
```

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
  val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
  w -= gradient
}</pre>
```



```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
  val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
  w -= gradient</pre>
```

}

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
   val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
   w -= gradient
}</pre>
```

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
  val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
  w -= gradient
}</pre>
```

val points = spark.textFile(...).

```
Computes a rtitioner(NODES)).
Gradient per
Var w = Vector.random(data point
for (i <- 1 to ITERATIONS) {
  val gradient = points.map(computeGradient(_,w)).reduce(_ + _)</pre>
```

```
w -= gradient
```

val points = spark.textFile(...).

```
wap(parsePoint).
p.Computes a Partitioner(NODES)).
Gradient per
data point
for (i <- 1 to ITERATIONS) {
 val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
 w -= gradient
```

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
   val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
   w -= gradient
}</pre>
```

```
val points = spark.textFile(...).
    map(parsePoint).
    partitionBy(HashPartitioner(NODES)).
    cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {
   val gradient = points.map(computeGradient(_,w)).reduce(_ + _)
   w -= gradient</pre>
```

}

```
val points = spark.textFile(...).
                   map(parsePoint).
                   partitionBy(HashPartitioner(NODES)).
                   cache()
```

```
var w = Vector.random(D)
```

```
for (i <- 1 to ITERATIONS) {</pre>
 val gradient = points.map(computeGradient(_,w))
```

w -= gradient

}

Trouble! shows through)

Spark: Conclusion

• The Good

- Speed (ca. MPI speed)
- Fault Tolerance
- Ease of Programming
- The Bad
 - Main Memory Assumption
- The Ugly
 - Systems aspects creep up



- Graph Computation framework
- Developed by Google
- Per vertex function

 update() processes
 incoming messages and
 sends new ones
- Computation is Bulk Synchronous Parallel

Pregel: A System for Large-Scale Graph Processing

Grzegorz Malewicz, Matthew H. Austern, Aart J. C. Bik, James C. Dehnert, Ilan Horn, Naty Leiser, and Grzegorz Czajkowski Google, Inc. {malewicz,austern,ajcbik,dehnert,ilan,naty,gczaj}@google.com

ГКАСТ

practical computing problems concern large graphs. ırd examples include the Web graph and various sotworks. The scale of these graphs—in some cases bilf vertices, trillions of edges-poses challenges to their it processing. In this paper we present a computamodel suitable for this task. Programs are expressed equence of iterations, in each of which a vertex can messages sent in the previous iteration, send meso other vertices, and modify its own state and that of going edges or mutate graph topology. This vertexapproach is flexible enough to express a broad set of hms. The model has been designed for efficient, scalid fault-tolerant implementation on clusters of thouof commodity computers, and its implied synchronickes reasoning about programs easier. Distributionl details are hidden behind an abstract API. The result mework for processing large graphs that is expressive sy to program.

gories and Subject Descriptors

Programming Techniques]: Concurrent Program-*Distributed programming*; D.2.13 [Software Enging]: Reusable Software—*Reusable libraries*

ral Terms

, Algorithms

/ords

outed computing, graph algorithms

NTRODUCTION

Internet made the Web graph a popular object of is and research. Web 2.0 fueled interest in social net-Other large graphs—for example induced by transion routes, similarity of newspaper articles, paths of

ion to make digital or hard conies of all or part of this work for

disease outbreaks, or citation relationships among public scientific work—have been processed for decades. Frequing applied algorithms include shortest paths computation ferent flavors of clustering, and variations on the page theme. There are many other graph computing proof practical value, *e.g.*, minimum cut and connected connects.

Efficient processing of large graphs is challenging. (algorithms often exhibit poor locality of memory access little work per vertex, and a changing degree of paral over the course of execution [31, 39]. Distribution over machines exacerbates the locality issue, and increase probability that a machine will fail during computation spite the ubiquity of large graphs and their commerciportance, we know of no scalable general-purpose sy for implementing arbitrary graph algorithms over arb graph representations in a large-scale distributed enment.

Implementing an algorithm to process a large graph ically means choosing among the following options:

- 1. Crafting a custom distributed infrastructure, typ requiring a substantial implementation effort that be repeated for each new algorithm or graph rep tation.
- 2. Relying on an existing distributed computing plat often ill-suited for graph processing. MapReduce for example, is a very good fit for a wide array of scale computing problems. It is sometimes us mine large graphs [11, 30], but this can lead tc optimal performance and usability issues. The models for processing data have been extended cilitate aggregation [41] and SQL-like queries [4(but these extensions are usually not ideal for gragorithms that often better fit a message passing n
- Using a single-computer graph algorithm library as BGL [43], LEDA [35], NetworkX [25], JDSI Stanford GraphBase [29], or FGL [16], limitin scale of problems that can be addressed.
- 4. Using an existing parallel graph system. The Pa BCL [22] and CCMgraph [8] librarios address pr

Giraph

- Apache Open Source implementation of Pregel
- Runs on Hadoop, (ab)uses mappers to do so



 Used at LinkedIn and Facebook

http://incubator.apache.org/giraph/

Messages Arrive and Are Processed

t=0

Messages Arrive and Are Processed






















Pregel: PageRank

- update() receives the PageRank of all neighbors
- Updates its local PageRank
- Sends new PageRank around if it changed enough

Pregel: Conclusion

- The Good
 - Excellent Map for Graph problems
 - Fast
- The Bad
 - Memory Model
 - Main Memory Assumption
- The Ugly
 - Wrong computational model (stay for the afternoon)

Open Problems

- No complete isolation of user / systems code
 - Unlike MapReduce
- No one system for example formation and modeling
 - Learning Effort
 - Orchestration
 - Wasted resources in distributed clusters



Joint Work With





Yingyi Bu, Vinayak Borkar, Michael J. Carey University of California, Irvine

Joshua Rosen, Neoklis Polyzotis University of California, Santa Cruz



Joshua Rosen, Neoklis Polyzotis University of California, Santa Cruz



- Unify Example Formation and Modeling
 - Relational Algebra Operators
 - Iteration Support
 - A unified runtime
- Increase Productivity via high-level language
 - Insulate the user from the systems aspects
 - Debugging and IDE support





High Level Language Relational Algebra and Loops



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime Scalability + High performance



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime Scalability + High performance

ScalOps

High Level Language Relational Algebra and Loops

Datalog Recursive Dataflow Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime Scalability + High performance

ScalOps

- Internal Domain Specific Language (DSL)
 - Written in Scala
- Relational Algebra (Filter, Join, GroupBy, ...)
- Iteration support
- Rich UDF support
 - Inline Scala function calls / literals
 - Byte-code compatible with Java
- Support in major IDEs

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
    val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
    val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
    env.w -= gradient
    env.delta = env.lastLoss - loss
    env.lastLoss = loss
    env
  }
}
```

def train(xy:Table[Example],

compute_grad:(Example, Vector) => Vector, compute_loss:(Example, Vector) => Double) = {

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
  }
}
```



class Example(x:Vector, y:Double)



```
def train(xy:Table[Example],
```

```
compute_grad:(Example, Vector) => Vector,
compute_loss:(Example, Vector) => Double) = {
```

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
  }
}
```



class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
 }
}
```



class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
    val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
    val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
    env.w -= gradient
    env.delta = env.lastLoss - loss
    env.lastLoss = loss
    env
  }
}
```







class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
}
```







class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
  }
}
```
Shared Loop State

ompute_loss:(Example, Vector) => Double) = {

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

le, Vector) => Vector,

val initialValue = new Env(VectorType.zeros(1000), Double.MaxValue, Double.MaxValue)

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
    val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
    val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
    env.w -= gradient
    env.delta = env.lastLoss - loss
    env.lastLoss = loss
    env
  }
}
```

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

val initialValue = new Env(VectorType.zeros(1000), Double.MaxValue, Double.MaxValue)

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
  }
}
```







Approach



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime Scalability + High performance

Approach



class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

val initialValue = new Env(VectorType.zeros(1000), Double.MaxValue, Double.MaxValue)

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
    val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
    val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
    env.w -= gradient
    env.delta = env.lastLoss - loss
    env.lastLoss = loss
    env
  }
}
```

Automatic Optimizations

class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

val initialValue = new Env(VectorType.zeros(1000), Double.MaxValue, Double.MaxValue)

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
}
```

Automatic Optimizations



class Env(w:VectorType, lastError:DoubleType, delta:DoubleType) extends Environment

val initialValue = new Env(VectorType.zeros(1000), Double.MaxValue, Double.MaxValue)

```
loop(initialValue, (env: Env) => env.delta < eps) { env => {
   val gradient = xy.map(x=>compute_grad(x,env.w)).reduce(_+_)
   val loss = xy.map(x=>compute_loss(x,env.w)).reduce(_+_)
   env.w -= gradient
   env.delta = env.lastLoss - loss
   env.lastLoss = loss
   env
}
```



Logical Plan



Approach



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime Scalability + High performance

Approach



High Level Language Relational Algebra and Loops

Declarative Language Captures the Recursive Dataflow

Suite of data-parallel operators Selected by an Optimizer / Compiler

Unified Runtime

Scalability + High performance

Monday, April 16, 2012

Some Optimizations

- Caching, Rocking
- Scheduling: Data-Local, Iteration-Aware
- Avoid (de-)serialization
- Minimize #network connections
- Pipelining

Physical Plan



Iterative Computation





Iterative Computation





Fan-In



Fan-In



Fan-In: Blocking



Fan-In: Blocking

$$h = \log_f(N) = \frac{\ln(N)}{\ln(f)}$$



Fan-In: Time per Level

$$h = \log_f(N) = \frac{\ln(N)}{\ln(f)}$$



Fan-In: Time per Level

$$t = fA$$

$$h = \log_f(N) = \frac{\ln(N)}{\ln(f)}$$

$$+$$

$$+$$

$$+$$

$$+$$

$$+$$

$$+$$

$$+$$

$$+$$

Fan-In: Total Time



Fan-In: Total Time

$$h = \log_f(N) = \frac{\ln(N)}{\ln(f)}$$



t = fA

Fan-In: Total Time



Partitioning

- Aggregation time increases logarithmically with number of machines
- Map time decreases linearly with the number of machines
- Closed form solutions available (but omitted here)

- As fast as
 - Vowpal Wabbit
 - Spark
- Faster than Hadoop (doh!)
- Much, much less code







Summary

- Example Formation
 - Use Pig
- Modeling
 - Hadoop (maybe not)
 - Subsampling (now)
 - Spark / Pregel (now)
 - ScalOps (as soon as we are done)