# cs 260D – Large Scale ML

#### Lecture #3: Distributed SGD

Baharan Mirzasoleiman UCLA Computer Science

Adapted from Dimitris Papailiopoulos

# Advances and Challenges in Distributed Machine Learning

Adapted from Dimitris Papailiopoulos

# Mostly Challenges



- Multicore vs. Distributed
- Algorithms of choice
- Open challenges with Performance Gains/Analysis
- Convergence/Generalization/Speedup/Delays?
- Straggler Nodes
- Adversarial attacks

# Stochastic Gradient Descent

Different names and flavors

ML / Optimization / Statistics / EE

Perceptron Incremental Gradient Back Propagation (NNs) Oja's iteration (PCA) LMS Filter

Has been around for a while, for good reasons: - Robust to noise - Simple to implement - Near-optimal learning performance \* - Small computational foot-print

. . .

# Stochastic Gradient Descent

#### Months! SGD can take 104 days on large data sets [DawnBench, 2017]

# Goal: Speed up Machine Learning

# Idea: Train at scale



### Platforms



# Parallel vs. Distributed

- Parallel (CPU)
  - Single machine, many cores (usually up to 100s)
  - Shared memory (all cores have access to RAM)
  - Comm. to RAM is cheap

Communication

- Distributed
  - Many machines (usually up to 1000s) connected via network
  - Shared-nothing architecture (each node has its own resources)
  - Communication costs non-negligible

### Scaling up vs. Scaling out

# Scaling up vs. out

#### What we'd ideally like

- Cores
- $\infty \\ \infty$ • RAM
- Comm. Cost 0
- 0 Cost to build

#### Feasible solutions:

Scaling up:

Getting the largest machine possible, with maxed out RAM

Scaling out:

Getting a bunch of machines, and linking them together

# Scaling up vs. out

#### Scaling up

#### <u>Pros:</u> RAM comm. cheap (no network) Less impl. overheads Less power/smaller footprint

#### <u>Cons:</u> 100s cores/machine = expensive f one fails, the system is dead Smaller fault tolerance Limited upgradability

#### Scaling out

#### <u>Pros:</u>

Much cheaper (especially on Ec2) Can replace faulty parts Better fault tolerance (if it matters)

#### <u>Cons:</u>

Network bound Major implementation overheads Large power footprint

# Choosing your Hardware

### Price Comparisons for 4 GPUs racks

- Single machine multi-core (CPU)
   For 4 Tesla V100 GPU work station: A100
   More than \$60,000
- Alternative Choice: Rent instances



**Google** Cloud Platform



#### EC2 Price for Scaling-up to 128 cores

• Alternative Choice: Rent Instances

Instance Size	GPUs - Testa V100 -		GPU Memory (GB)	vcpus	Memory (GB)	Network Bandwidth	EBS Bandwidth	On-Demand Price/hr*	1-yr Reserved Instance Effective Hourly*	3-yr Reserved Instance Effective Hourly*
p3.8xlarge	4	NVLink	64	32	244	10 Gbps	7 Gbps	\$12.24	\$7.96	\$4.93
p3.16xlarge	8	NVLink	128	64	488	25 Gbps	14 Gbps	\$24.48	\$15.91	\$9.87

Price for 4 GPUs: \$8/hr

10 days nonstop/month, for a year ~ \$23,000

# Let's Think about Running Stuff

How to Parallelize  

$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k}(w_k; x_i)$$

- **Parallelize** computation of **one update**? <u>Issue:</u> computing  $\nabla \ell_i(w; x_i)$  cheap (even for deep nets) [O(d)]
  - Parallel Updates? 所谓串行,指的是一个mini-batch完成后, 再进行下一个mini-batch Issue: SGD is inherently serial....

Q: Can we parallelize inherently serial algorithms?

### Simple idea

Compute multiple gradients in parallel

K应该是epoch的计数

$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^1}(w_k; x_i)$$
$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^2}(w_k; x_i)$$
$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^3}(w_k; x_i)$$
$$w_{k+1} = w_k - \gamma \nabla \ell_{s_k^4}(w_k; x_i)$$

#### Q: Does it perform the same as SGD?

# Distribute the effort!













#### 每个机器上都有一个小模型

- Many models weaker than one
- Delays and Slow Nodes Theory

- **Communication Costs**
- Barriers to entry / High Cost
- Implementation Overhead Practice
- Nontrivial choice of ML-framework

#### Several issues

# Choosing an Algorithm

# The Master-worker Setting















Repeat distributed iterations until we are happy with the model . . . .

# Evaluating the performance of Distributed ML

### Many Questions....

- How fast does minibatch-SGD converge?
- How can we measure speedups? Average指的是在master node上进行合并权重
- Comm. is expensive, how often do we average?
- How do we choose the right model?
- What happens with delayed nodes?
- Does fault tolerance matter?



# Why so slow?

#### Two factors control run-time

#### Time to accuracy $\varepsilon =$ [time per data pass] X [#passes to accuracy $\varepsilon$ ]

Time cost of each epoch

**Epoch numbers** 

# Why so slow?

"...] on more than 8 machines [...] network overhead starts to dominate [...]

TL;DR: Communication is the Bottleneck Larger batch size faster but more expensive on communication
 Why?



Why is this better than optimum

More worker, less batch numbers

Time per pass: time for dataset\_size/batch\_size distributed iterations

> Bigger Batch ⇒ Less Communication (smaller time per epoch)

## What's wrong with Large Batches?

- If small batch is bad, then maximize it



# How to analyze parallel algorithms? Main measure of performance Time of serial $\mathcal{A}$ to accuraccy $\epsilon$ speedup = $\frac{1}{\text{Time of parallel } \mathcal{A} \text{ to accuraccy } \epsilon}$ Example: Gradient Descent Embarrassingly Parallel O(#cores) speedup

- Convergence is **invariant of allocation** 

Both algorithNot true for mini-batch SGD terT iterations
 Speedup is independent of covergence rate

# How to Analyze mini-batch?

• Measure of performance Epsilon is the accuracy goal

worst case speedup =  $\frac{\text{bound on \#iter of SGD to }\epsilon}{\text{bound on #iter of Parallel SGD to }\epsilon}$ 

Main Question: How does minibatch SGD compare against serial SGD?

Main questions:

Convergence after T gradient computations

Answer is Complicated: Depends on Problem Generally if batch B\_0 > B(data, loss) Minibatch SGD offers no speedup.

Diversity related

#### Generalization?

#### Do models trained with minibatch SGD generalize

#### You can't train your neural network on too large batch size

Batch	Top-1 Acc	Top-5 Acc
256	58.42%	81.51%
512	59.19%	81.84%
1024	59.00%	81.94%
2048	58.88%	81.73%
4096	57.97%	81.00%
8192	55.90%	79.40%

Alexnet on Imagenet
#### How to choose the right model

Some models are better than others, even from a systems perspective

#### How to choose the right model

Some models are better than others, even from a systems perspective



It's fully connected so it's hard to find a place to split

被切断的链接仍然存在,并没有消失(具体实现 不清楚)

#### How to choose the right model

Some models are better than others, even from a systems perspective



- Does it fit in a single machine?
- Is model architecture amenable to low communication?
- Some models easier to partition
- Can we increase sparsity (less comm) without losing with accuracy?

#### How to chose the right model

Some models are better than others, even from a systems perspective



#### Parallelism Gains?

• Weak Scaling: Easy

• Strong Scaling: Nontrivial, and not there yet

#### Major Open Problem!!!

# Avoiding Communication Bottlenecks

#### Large-scale Distributed Machine Learning Systems



## A better algorithm?



# **QSGD:** Communication-Efficient SGD via Gradient Quantization and Encoding





Parallel Stochastic Gradient Descent (PSGD)

Shuffle data in different machines cost a lot

Data shuffling

# Very high communication cost: does shuffling still makes sense?



\* [Recht and Re, 2013], [Bottou, 2012], [Zhang and Re, 2014], [Gurbuzbalaban et al., 2015], [loffe and Szegedy, 2015], [Zhang et al. 2015]

# Bottleneck 2: Straggling Learners



# A case against Synchronization



#### Asynchronous World



#### Stragglers

- Ideal compute time per node ~ O(total\_time/P)
- But there is a lot of randomness:
  - Network/Comm Delays
  - Node/HW Failures
  - Resource Sharing
- What if time per node is a random variable:  $X = constant + Exp(\lambda)$



#### Remark:

Slowest node is log(n) times slower than fastest

#### Simulation

• X(t) = I + Exp(0.5), n = I0, I00, I000, I000



# Parallelizing Sparse SGD on shared memory architectures

#### Single Machine, Multi-core



# Asynchronous SGD on Sparse Functions

$$f(x) = \sum_{e \in \mathcal{E}} f_e(x_e)$$

- Def: Hyperedge e = the subset of variables that  $f_e$  depends on
- The function-variable graph







Pick random data point



Read Variables



Step3: Compute local function



# Examples of Sparse Problems

### Sparse Support Vector Machines



#### Matrix Completion



Entries Specified on set *E* (with |E|=n)

$$\operatorname{minimize}_{(\mathbf{L},\mathbf{R})} \sum_{(u,v)\in E} \left\{ (\mathbf{L}_u \mathbf{R}_v^T - M_{uv})^2 + \mu_u \|\mathbf{L}_u\|_F^2 + \mu_v \|\mathbf{R}_v\|_F^2 \right\}$$

## Graph Cuts



- Image Segmentation
- Entity Resolution
- Topic Modeling

minimize<sub>x</sub>  $\sum_{(u,v)\in E} w_{uv} ||x_u - x_v||_1$ subject to  $\mathbf{1}_K^T x_v = 1, x_v \ge 0, \text{ for } v = 1, \dots, D$ 

### Sparsified BackProp

Dropout: drop some nodes



<sup>(</sup>a) Standard Neural Net



(b) After applying dropout.

#### Challenges in Parallel SGD



No conflict =>2 parallel iterations = 2 serial iterations

#### Challenges in Parallel SGD



#### No conflict => Speedup

#### Challenges in Parallel SGD Shared RAM



What should we do for conflicts? <u>Approach I</u>: Coordinate or Lock <u>Approach 2</u>: Don't Care (Lock-free Async.)

# Prior to 2011 Work

Long line of theoretical work since the 60s [Chazan, Miranker, 1969]

Foundational work on Asynchronous Optimization Master/Worker model [Tsitsiklis, Bertsekas, 1986, 1989]

Recent hardware/software advances renewed the interest Round-robin approach [Zinkevich, Langford, Smola, 2009] Average Runs [Zinkevich et al., 2009], Average Gradients [Duchi et al, Dekel et al. 2010]



#### HOGWILD! 2011 "('Run parallel lock-free SGD without synchronization"













Google Downpour SGD, **Microsoft** Project Adam use HOGWILD! Renewed interest on async. optimization

#### Speedups



Experiments run on 12 core machine 10 cores for gradients, 1 core for data shuffling

# Challenges in Analysis

## Challenges in Hogwild!

Shared Memory



#### Incompatible with classic SGD analysis
## Analyzing Asynchronous Schemes

## A Noisy Lens <sup>Noice的视角</sup> for Asynchronous Algorithms

### Main Idea

#### Noisy viewpoint: Asynchronous(Algo.(INPUT )) $\equiv$ Serial(Algo.(INPUT + Noise)

#### 扰乱的

Perturbed Iterate Analysis for Asynchronous Stochastic Optimization [Mania, Pan, P, Recht, Ramchandran, Jordan, 2015]

Joint work with











Each processor in parallel sample function  $f_i$ x = read shared memory $g = -\gamma \cdot \nabla f_i(x)$ for v in the support of f do  $x_v \leftarrow x_v + g_v$ 

- Def:  $S_k$  is the k-th sampled data point
- Fact: Cores don't read "actual" iterates  $x_k$  but "noisy iterates"  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)

Ex.



Each processor in parallel sample function  $f_i$ x = read shared memory $g = -\gamma \cdot \nabla f_i(x)$ for v in the support of f do  $x_v \leftarrow x_v + g_v$ 

- Def:  $S_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates  $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)





Each processor in parallel sample function  $f_i$ x = read shared memory $g = -\gamma \cdot \nabla f_i(x)$ for v in the support of f do  $x_v \leftarrow x_v + g_v$ 

- Def:  $S_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After *T* processed samples, the contents of RAM are: *(atomic writes + commutativity)*







- <u>Def:</u>  $s_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates  $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)





Each processor in parallel sample function  $f_i$ x = read shared memory $g = -\gamma \cdot \nabla f_i(x)$ for v in the support of f do  $x_v \leftarrow x_v + g_v$ 

- Def:  $S_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)







- <u>Def:</u>  $s_k$  is the k-th sampled data point
- Fact: Cores don't read ''actual'' iterates  $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)







- <u>Def:</u>  $S_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)





Each processor in parallel sample function  $f_i$ x = read shared memory $g = -\gamma \cdot \nabla f_i(x)$ for v in the support of f do  $x_v \leftarrow x_v + g_v$ 

- Def:  $S_k$  is the k-th sampled data point
- Fact: Cores don't read ''actual'' iterates  $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After *T* processed samples, the contents of RAM are: *(atomic writes + commutativity)*







- <u>Def:</u>  $s_k$  is the k-th sampled data point
- <u>Fact:</u> Cores don't read ''actual'' iterates  $x_k$  but ''noisy iterates''  $\hat{x}_k$
- After T processed samples, the contents of RAM are: (atomic writes + commutativity)

$$x_0 - \gamma \cdot \nabla f_{s_0}(\hat{x}_0) - \ldots - \gamma \cdot \nabla f_{s_{T-1}}(\hat{x}_{T-1})$$

Main Questions:1) Where does noise come from?2) How strong is it?







"Serialized" Processing Timeline







#### "Serialized" Processing Timeline







#### "Serialized" Processing Timeline







#### "Serialized" Processing Timeline



Asynchrony noise is <u>combinatorial</u> coordinates in conflict can be as noisy as possible. (no generative model assumptions)



### Convergence of Hogwild

THEOREM 3.4. If the number of samples that overlap in time with a single sample during the execution of HOGWILD! is bounded as

$$\tau = \mathcal{O}\left(\min\left\{\frac{n}{\overline{\Delta}_C}, \frac{M^2}{\epsilon m^2}\right\}\right),$$

HOGWILD!, with step size  $\gamma = \frac{\epsilon m}{2M^2}$ , reaches an accuracy of  $\mathbb{E} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \epsilon$  after

$$T \ge \mathcal{O}(1) \frac{M^2 \log\left(\frac{a_0}{\epsilon}\right)}{\epsilon m^2}$$

iterations.

## Open Problems

### Open Problems: Part I

Assumptions:

Sparsity + convexity => linear speedups

O.P. : Hogwild! On Dense Problems

> Maybe we should featurize dense ML Problems, so that updates are **sparse**

> > O.P.:

Fundamental Trade-off Sparsity vs Learning Quality?

### Open Problems: Part 2

- What we proved:

worst case speedup =  $\frac{\text{bound on \#iter of SGD to }\epsilon}{\text{bound on \#iter of Parallel SGD to }\epsilon}$ 

- What we really care about:

speedup =  $\frac{\text{Time of serial } \mathcal{A} \text{ to accuraccy } \epsilon}{\text{Time of parallel } \mathcal{A} \text{ to accuraccy } \epsilon}$ 

O.P. : True Speedup Proofs for Hogwild

Holy Grail

O.P. : Guarantees for Nonconvex Problems?

### Open Problems: Part 3

Hogwild! Algorithms great for Shared Memory Systems



- Similar Issues for Distributed:

O.P. : Sync vs Async is still open

## Reproducible Models

# Reproducibility

因为结束时间随机,覆盖的顺序不一样

- HOGWILD! Models are not reproducible
- Each training session has inherent "system" randomness
- Does not allow to recreate models if needed
- 可解释性 - Barrier for accountability and reproducibility
- How can we resolve it?

## Presentation Suggestion

Suggestion: Minibatch, Hogwils Analysis? Hadoop?

## Next Week's Presentation

- Polyak, Boris T. "Some methods of speeding up the convergence of iteration methods." *Ussr computational mathematics and mathematical physics* 4.5 (1964): 1-17.
- Sutskever, Ilya, et al. "On the importance of initialization and momentum in deep learning." *International conference on machine learning*. PMLR, 2013.

Why Momentum Really Works, Distill, [Goh, 2017]. <u>http://doi.org/10.23915/distill.00006</u>

Martens, James, and Roger Grosse. "Optimizing neural networks with kroneckerfactored approximate curvature." *International conference on machine learning*. PMLR, 2015.