

## Case Study 3: fMRI Prediction

“Scalable” LASSO Solvers:  
Parallel SCD (Shotgun)  
Parallel SGD  
Averaging Solutions  
ADMM

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

February 6<sup>th</sup>, 2014

©Emily Fox 2014

1

## Scaling Up LASSO Solvers

- A simple SCD for LASSO (Shooting)
  - Your HW, a more efficient implementation! ☺
  - Analysis of SCD
- Parallel SCD (Shotgun)
- Other parallel learning approaches for linear models
  - Parallel stochastic gradient descent (SGD)
  - Parallel independent solutions then averaging
- ADMM

©Emily Fox 2014

2

# Stochastic Coordinate Descent for LASSO (aka Shooting Algorithm)

## Repeat until convergence

- Pick a coordinate  $j$  at random

Set: 
$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases} = \text{sign}(c_j) \left( \frac{|c_j| - \lambda}{a_j} \right)$$

Where: 
$$a_j = 2 \sum_{i=1}^N (x_j^i)^2 \quad c_j = 2 \sum_{i=1}^N x_j^i (y^i - \beta'_{-j} x_{-j}^i)$$

cost per iteration  $O(N)$

Can be done more efficiently. Proof: HW!

$\min_{\beta} F(\beta_1, \dots, \beta_{j-1}, \beta, \beta_{j+1}, \dots, \beta_p)$   
 $\uparrow$   
 $j^{\text{th}}$  coord.

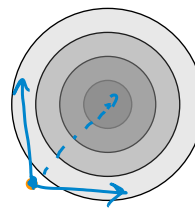
# Shotgun: Parallel SCD [Bradley et al '11]

Lasso:  $\min_{\beta} F(\beta)$  where  $F(\beta) = \|X\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1$

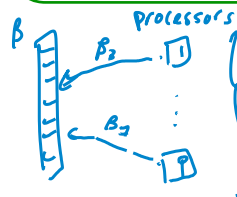
## Shotgun (Parallel SCD)

While not converged,

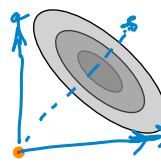
- On each of  $P$  processors,
- Choose random coordinate  $j$ ,
- Update  $\beta_j$  (same as for Shooting)



yes!  
Features are uncorrelated



act as if they were the only children

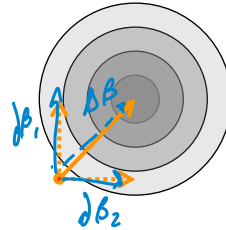


no!!  
Features are highly corr.

# Is SCD inherently sequential?

Lasso:  $\min_{\beta} F(\beta)$  where  $F(\beta) = \|X\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1$

Coordinate update:  
 $\beta_j \leftarrow \beta_j + \delta\beta_j$   
 (closed-form minimization)



Collective update:

$$\Delta\beta = \begin{pmatrix} \delta\beta_i \\ 0 \\ 0 \\ \delta\beta_j \\ 0 \end{pmatrix}$$

there are interferences in these updates if features are corr.  
 Can we quantify this?

©Emily Fox 2014

5

# Convergence Analysis

Lasso:  $\min_{\beta} F(\beta)$  where  $F(\beta) = \|X\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1$

Theorem: Shotgun Convergence

Assume  $P < \frac{p}{\rho + 1}$

where  $\rho =$  spectral radius of  $X^T X$

$E[F(\beta^{(T)})] - F(\beta^*)$

$\leq \frac{p \left( \frac{1}{2} \|\beta^*\|_2^2 + F(\beta^{(0)}) \right)}{TP}$

↑  
 speed up linear in # proc,  
 up to  $P_{max}$

Nice case:  
 Uncorrelated features



$\rho = 1 \Rightarrow P_{max} = p$

Bad case:  
 Correlated features



$\rho = p \Rightarrow P_{max} = 1$  (at worst)

©Emily Fox 2014

6

# Stepping Back...

- Stochastic coordinate ascent <sup>SCD</sup>
  - Optimization: pick a coord.  $j$ , find  $\min_{\beta_j}$
  - Parallel SCD: pick  $P$  coord.
  - Issue: coordinates may interfere on  $P$  coord.  $\swarrow$  spectral radius
  - Solution: bound possible interference based  $\rho$
- Natural counterpart: SGD
  - Optimization: pick a datapoint  $i$   $\beta \leftarrow \beta - \eta \nabla F(x^i; \beta)$
  - Parallel: pick  $P$  datapoints + ind. update  $\beta$
  - Issue: can interfere on all coord.
  - Solution: bound interference by exploiting sparsity in  $X$

©Emily Fox 2014

7

# Parallel SGD with No Locks

[e.g., Hogwild!, Niu et al. '11]

- Each processor in parallel:
  - Pick data point  $i$  at random
  - For  $j = 1 \dots p$ :

$$\beta_j \leftarrow \beta_j - \eta (\nabla F(x^i; \beta))_j$$

- Assume atomicity of:  $\beta_j \leftarrow \beta_j + a$   
other interferences

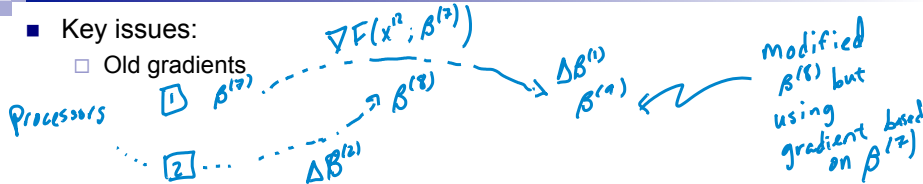
©Emily Fox 2014

8

# Addressing Interference in Parallel SGD

- Key issues:

- Old gradients



- Processors overwrite each other's work

- Nonetheless:

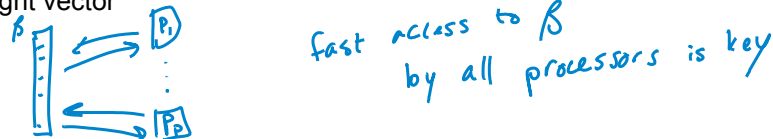
- Can achieve convergence and parallel speedups
- Proof uses weak interactions, but through sparsity of data points

Sparsity of  $X$  is key to the analysis

Update is exact for two  $x$ 's that do not share any support pts.

# Problem with Parallel SCD and SGD

- Both Parallel SCD & SGD assume access to current estimate of weight vector



fast access to  $\beta$  by all processors is key

- Works well on shared memory machines multicore

- Very difficult to implement efficiently in distributed memory cloud



- Open problem: Good parallel SGD and SCD for distributed setting...

- Let's look at a trivial approach

some work very recently

# Simplest Distributed Optimization Algorithm Ever Made

- Given  $N$  data points &  $P$  machines
- Stochastic optimization problem:
- Distribute data:  $P$  machines

$$\min_{\beta} F(\beta) \equiv \frac{1}{N} \sum_{i=1}^N F(x^i; \beta)$$

randomly  
assign data



solves a problem  $D_k$

$$|D_k| = \frac{N}{P} = n$$

- Solve problems independently

machine  $k$  : ind. est.  $\beta^{(k)} = \min_{\beta} \frac{1}{n} \sum_{i \in D_k} F(x^i; \beta)$

- Merge solutions

$$\bar{\beta} = \frac{1}{P} \sum_k \beta^{(k)}$$

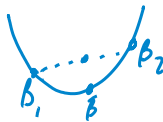
- Why should this work at all????

©Emily Fox 2014

11

# For Convex Functions...

- Convexity:



$$\frac{F(\beta_1) + F(\beta_2)}{2} \geq F(\bar{\beta})$$

- Thus:

$$\max(F(\beta_1), F(\beta_2)) \geq F(\bar{\beta})$$

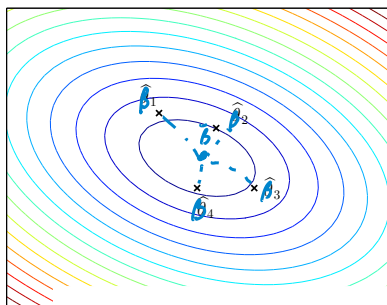


©Emily Fox 2014

12

# Hopefully...

using convexity alone



- Convexity only guarantees:

$$F(\bar{\beta}) \leq \max_k F(\beta^{(k)})$$

- But, estimates from independent data!

can we leverage this to improve this bound?

Figure from John Duchi

# Analysis of Distribute-then-Average

[Zhang et al. '12]

- Under some conditions, including strong convexity, lots of smoothness, and more...  
 $\hat{\beta}_N \in \arg \min_{\beta} \frac{1}{N} \sum_{i=1}^N F(x^i; \beta)$

- If all data were in one machine, converge at rate:

$$E[\|\hat{\beta}_N - \beta^*\|_2^2] = O\left(\frac{1}{N}\right)$$

- With  $P$  machines, converge at a rate:

$$E[\|\bar{\beta} - \beta^*\|_2^2] = O\left(\frac{1}{N} + \frac{1}{n^2}\right)$$

unavoidable

"bias" from parallelism

\* obs. per machine  
 $n = \frac{N}{P}$  ← # obs.  
 $P$  ← # proc.

e.g. 1T datapoints, 1000 machines  $\Rightarrow P = N^{1/4}$   
 Plug in  $\frac{1}{n^2} = \frac{1}{(N^{1/4})^2}$  ← negligible compared to  $\frac{1}{N}$ ...  
 great parallelism

# Tradeoffs, tradeoffs, tradeoffs,...

- Distribute-then-Average:

- "Minimum possible" communication
- Bias term can be a killer with finite data
  - Issue definitely observed in practice
- Significant issues for L1 problems:

*all ind. problems on each machine ... just merge at end*  
*prev. results are asy.*  
*sparsity patterns in machine i can be very diff. from those in machine j ⇒ average β can lose sparsity*

- Parallel SCD or SGD

- Can have much better convergence in practice for multicore setting
- Preserves sparsity (especially SCD)
- But, hard to implement in distributed setting

# Alternating Directions Method of Multipliers

- A tool for solving convex problems with separable objectives:

$$\min_x \{ f(x) + g(x) \}$$

- LASSO example:

$$\min_{\beta} \left\{ \underbrace{\|y - X\beta\|_2^2}_{f(\beta)} + \lambda \underbrace{\|\beta\|_1}_{g(\beta)} \right\}$$

- Know how to minimize  $f(\beta)$  or  $g(\beta)$  separately  
*coupling presents challenges*

*ADMM = approach that works in a dist. setting, but requires more comm. than dist + avg. approach (but less than SGD)*



# ADMM Insight

- Try this instead:

$$\min_{x, z} \{ f(x) + g(z) \} \quad \text{s.t. } x = z$$

*still convex!*

- Solve using method of multipliers
- Define the augmented Lagrangian:

$$L_{\rho}(x, z, y) = f(x) + g(z) + y^T(x - z) + \frac{\rho}{2} \|x - z\|_2^2$$

*pos. const.*

- Issue: L2 penalty destroys separability of Lagrangian
- Solution: Replace minimization over  $(x, z)$  by alternating minimization

©Emily Fox 2014

17

# ADMM Algorithm

- Augmented Lagrangian:

$$L_{\rho}(x, z, y) = f(x) + g(z) + y^T(x - z) + \frac{\rho}{2} \|x - z\|_2^2$$

- Alternate between:

- $x \leftarrow \arg \min_x L_{\rho}(x, z, y)$
- $z \leftarrow \arg \min_z L_{\rho}(x, z, y)$
- $y \leftarrow y + \rho(x - z)$

©Emily Fox 2014

18

# ADMM for LASSO

$$L_\rho(x, z, y) = f(x) + g(z) + y^T(x - z) + \frac{\rho}{2}\|x - z\|_2^2$$

- Objective:  $\min_{\beta, z} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|z\|_1 \right\}$  s.t.  $\beta = z$

- Augmented Lagrangian:

$$L_\rho(\beta, z, a) = \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|z\|_1 + a^T(\beta - z) + \frac{\rho}{2} \|\beta - z\|_2^2$$

- Alternate between:

$$1. \beta \leftarrow \arg \min_{\beta} L_\rho(\beta, z, a) = \overbrace{(X^T X + \rho I)^{-1}}^{\text{precompute}} \overbrace{(X^T y + \rho z - a)}^{\text{distribute}}$$

$$2. z \leftarrow \arg \min_z L_\rho(\beta, z, a) = S\left(\beta + \frac{a}{\rho}, \frac{1}{\rho}\right)$$

$$3. a \leftarrow a + \rho(\beta - z)$$

$\leftarrow S(a, c) = \text{sign}(a) (|a| - c)_+$   
 soft-thresholding

For distributed version, see paper. .

©Emily Fox 2014

19

# ADMM Wrap-Up

- When does ADMM converge?

- Under very mild conditions
- Basically, f and g must be convex

- ADMM is useful in cases where

- $f(x) + g(x)$  is challenging to solve due to coupling
- We can minimize
  - $f(x) + (x-a)^2$
  - $g(x) + (x-a)^2$

- Reference

- Boyd, Parikh, Chu, Peleato, Eckstein (2011) "Distributed optimization and statistical learning via the alternating direction method of multipliers." *Foundations and Trends in Machine Learning*, 3(1):1-122.



©Emily Fox 2014

20

# What you need to know

- A simple SCD for LASSO (Shooting)
  - Your HW, a more efficient implementation! ☺
  - Analysis of SCD
- Parallel SCD (Shotgun)
- Other parallel learning approaches for linear models
  - Parallel stochastic gradient descent (SGD)
  - Parallel independent solutions then averaging
- ADMM
  - General idea
  - Application to LASSO