# Case Study 3: fMRI Prediction

**Stochastic Coordinate Descent (SCD) for LASSO (Shooting)**
**Parallel SCD (Shotgun)**
**Parallel SGD**
**Averaging Solutions**

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin

February 21st, 2013

1

---

# Today

- One way to solve LASSO problem
- Stochastic Coordinate Descent (SCD)
- Minimizing a coordinate in LASSO
- 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

2

# Coordinate Descent

$F(\beta_1, \ldots, \beta_\lambda)$

- Given a function $F(\beta)$
  - □ Want to find minimum   $\beta^* \leftarrow \min_{\beta} F(\beta)$

- Often, hard to find minimum for all coordinates, but easy for one coordinate

  1-d optimization problem

- Coordinate descent:

  while not converged
  Pick coordinate j
  $$\beta_j \leftarrow \min_{b} F(\beta_1, \beta_2, \ldots, \beta_{j-1}, b, \beta_{j+1}, \ldots, \beta_\lambda)$$

- How do we pick a coordinate?

  Round robin, random, smartly, ...

- When does this converge to optimum?   e.g., strongly convex (Separability)

**3**

---

# LASSO Regression

- **LASSO:** least absolute shrinkage and selection operator

- New objective:

$$\min_{\beta} \sum_{i=1}^{N} \left( y^i - (\beta_0 + \beta^T x^i) \right)^2 + \lambda \|\beta\|_1$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{RSS(\beta)}$$

$$\Updownarrow$$

$$\min_{\beta} RSS(\beta) \qquad s.t. \ \|\beta\|_1 \leq B$$

**4**

---

2

# Soft Threshholding

$F(\beta) = RSS(\beta) + \lambda \|\beta\|_1$

- Gradient of RSS term:

$$\frac{\partial}{\partial \beta_j} RSS(\beta) = a_j \beta_j - c_j$$

$$a_j = 2 \sum_{i=1}^{N} (x_j^i)^2$$

$$c_j \leftarrow 2 \sum_{i=1}^{N} x_j^i (y^i - \beta_{-j}^T x_{-j}^i)$$

all but the $j^{th}$ coeff.

$c_j \propto corr(x_j, r_{-j})$

msr how relevant $x_j$ is for pred $y$ beyond what the others can

residual from model w/o $j^{th}$ cov.

- Subgradient of full objective:

$$\partial_{\beta_j} F(\beta) = (a_j \beta_j - c_j) + \lambda \, \partial_{\beta_j} \|\beta\|_1$$

$$= \begin{cases} a_j \beta_j - c_j - \lambda & \beta_j < 0 \\ [-c_j - \lambda, -c_j + \lambda] & \beta_j = 0 \\ a_j \beta_j - c_j + \lambda & \beta_j > 0 \end{cases}$$

©Carlos Guestrin 2013

5

---

# Soft Threshholding

$\partial_{\beta_j} F = 0 \Rightarrow \min_{\beta_j} F(\beta_1, \ldots, \beta_{j-1}, \beta_j, \beta_{j+1}, \ldots, \beta_p)$

- Set subgradient = 0:

$$\partial_{\beta_j} F(\beta) = \begin{cases} a_j \beta_j - c_j - \lambda & \beta_j < 0 \\ [-c_j - \lambda, -c_j + \lambda] & \beta_j = 0 \\ a_j \beta_j - c_j + \lambda & \beta_j > 0 \end{cases}$$

$$= 0$$

If $\beta_j < 0$

$$a_j \beta_j - c_j - \lambda = 0$$

$$\Rightarrow \beta_j = \frac{c_j + \lambda}{a_j} < 0 \Rightarrow c_j < -\lambda$$

strong neg. corr., then $\beta_j < 0$

If $\beta_j > 0$

$$a_j \beta_j - c_j + \lambda = 0 \Rightarrow \beta_j = \frac{c_j - \lambda}{a_j} > 0 \Rightarrow c_j > \lambda$$

strong pos. corr., then $\beta_j > 0$

If $\beta_j = 0 \quad -\lambda < c_j < \lambda$ \qquad Otherwise, $\beta_j = 0$

- The value of $c_j = 2 \sum_{i=1}^{N} x_j^i (y^i - \beta'_{-j} x_{-j}^i)$ constrains $\beta_j$

©Carlos Guestrin 2013

6

---

3

# Soft Threshholding

$a_j \geq 0$

$\min_{\beta} f(\ldots \beta_{j-1}, \beta_j, \beta_{j+1}, \ldots)$
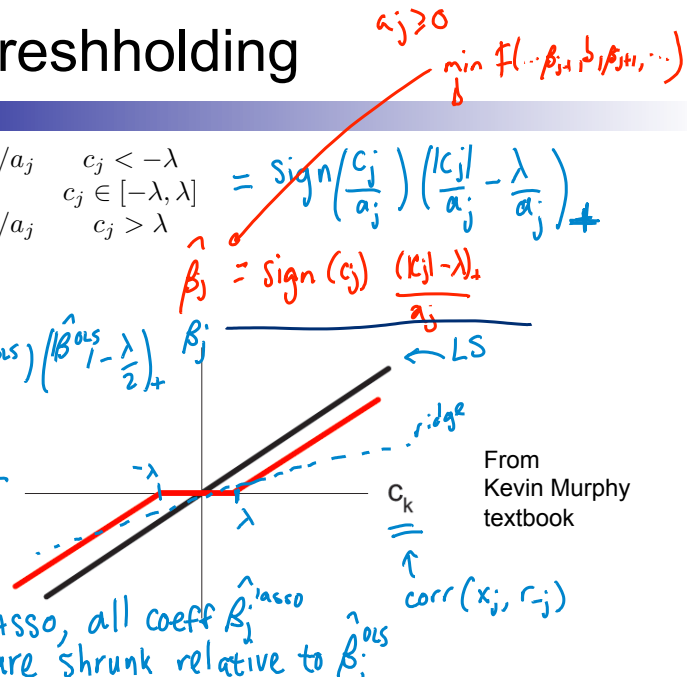
$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases} = \text{Sign}\left(\frac{c_j}{a_j}\right)\left(\frac{|c_j|}{a_j} - \frac{\lambda}{a_j}\right)_+$$

$\hat{\beta}_j = \text{Sign}(c_j) \dfrac{(|c_j| - \lambda)_+}{a_j}$

If $X^T X = I$

$\hat{\beta}_j^{lasso} = \text{Sign}\left(\hat{\beta}_j^{OLS}\right)\left(|\hat{\beta}_j^{OLS}| - \frac{\lambda}{2}\right)_+$

$\hat{\beta}_j^{ridge} = \dfrac{\hat{\beta}_j^{OLS}}{1+\lambda}$

← LS

ridge

$c_k$ = corr$(x_j, r_{-j})$

From Kevin Murphy textbook

In LASSO, all coeff $\hat{\beta}_j^{lasso}$ are shrunk relative to $\hat{\beta}_j^{OLS}$

7

---

# Stochastic Coordinate Descent for LASSO (aka Shooting Algorithm)

- Repeat until convergence
  - □ Pick a coordinate *j* at random
    - Set: $$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases} = \text{Sign}(c_j)\dfrac{(|c_j| - \lambda)_+}{a_j}$$
    - Where: 
      
      cache
      
      $$a_j = 2\sum_{i=1}^{N}(x_j^i)^2 \qquad c_j = 2\sum_{i=1}^{N}x_j^i(y^i - \beta'_{-j}x_{-j}^i)$$

cost per iteration

$O(N)$

proof: your HW!!

can be done more smartly ...

8

4

# Analysis of SCD [Shalev-Shwartz, Tewari '09/'11]

$e_j = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}$

- Analysis works for LASSO, L1 regularized logistic regression, and other objectives!

- For (coordinate-wise) strongly convex functions:

$$F(\beta + \Delta\beta) \leq F(\beta) + \partial_{\beta_j} (\nabla F(\beta))_j + \gamma (\partial_{\beta_j})^2$$

$\Delta\beta = \partial_{\beta_j} e_j$

dim   how hard   where $\pm$ start $^2$
         d

Lasso
$\gamma = 1$

Logistic
Regression
$\gamma = \frac{1}{4}$

- Theorem:      $\beta^{(0)}$
  - Starting from
  - After T iterations

dist from opt

$$E\left[F(\beta^{(T)})\right] - F(\beta^*) \leq \frac{d\left(\gamma \|\beta^*\|_2^2 + 2 F(\beta^{(0)})\right)}{T+1}$$

← gets linearly
better with iters

  - Where E[ ] is wrt random coordinate choices of SCD

- Natural question: How does SCD & SGD convergence rates differ?   no params to tune

See prev.:   SCD → faster w. larger d ← no params to tune
            SGD → faster w. larger N ← needs $\eta$

©Carlos Guestrin 2013          9

---

# Shooting: Sequential SCD

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

$F(\beta)$ contour

Stochastic Coordinate Descent (SCD)
(e.g., Shalev-Shwartz & Tewari, 2009)

While not converged,
- Choose random coordinate j,
- Update $\beta_j$ (closed-form minimization)

how do we measure?
→ annoying: over a time window?
                has anything changed?
↳ do a round robin iter to measure
                convergence

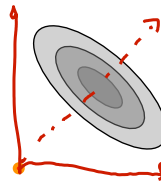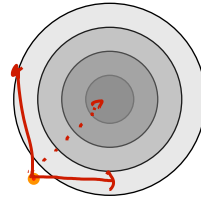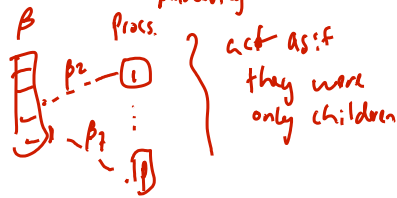©Carlos Guestrin 2013          10

5

# 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)

*independently*

$\beta$

Procs.

$\beta^2 \to \boxed{1}$

$\sim \beta^1 \to \boxed{P}$

act as if they were only children

yes!!

fectures are uncorrelated

no!!

fectures are highly correlted

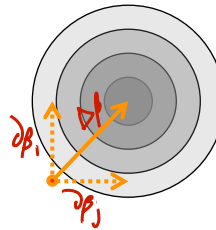©Carlos Guestrin 2013

11

---

# 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)

$\partial\beta_i$

$\partial\beta_j$

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

©Carlos Guestrin 2013

12

6

# Is SCD inherently sequential?

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

Theorem: If X is normalized s.t. diag($X^TX$)=1,

$F(\beta + \Delta\beta) - F(\beta)$ ⟶ *decrease in objective*

$$\leq -\sum_{i_j \in \mathcal{P}} \left(\delta\beta_{i_j}\right)^2 + \sum_{\substack{i_j, i_k \in \mathcal{P}, \\ j \neq k}} \left(X^TX\right)_{i_j, i_k} \delta\beta_{i_j} \delta\beta_{i_k}$$

"positive" progress

could be positive or negative

"interference" or "bias" from parallelism

13

---

# Is SCD inherently sequential?

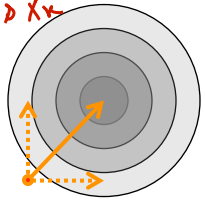Theorem: If X is normalized s.t. diag($X^TX$)=1,

$F(\beta + \Delta\beta) - F(\beta)$

Key term ⟵ measures magnitude of interference

$$\leq -\sum_{i_j \in \mathcal{P}} \left(\delta\beta_{i_j}\right)^2 + \sum_{\substack{i_j, i_k \in \mathcal{P}, \\ j \neq k}} \left(X^TX\right)_{i_j, i_k} \delta\beta_{i_j} \delta\beta_{i_k}$$
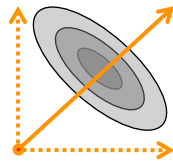
$\left(X^TX\right)_{jk} = 0$

"correlation" $x_j$ & $x_k$

$\left(X^TX\right)_{jk} \neq 0$

"interference"

Nice case: Uncorrelated features

Bad case: Correlated features

14

7

# Shotgun: Convergence Analysis

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

Assume # parallel updates $P < d/\rho + 1$

*dim*

*Spectral radius of $X^T X$*

$$E\left[F(\beta^{(t)})\right] - F(\beta^*) \leq \frac{d\left(\|\beta^*\|_2^2 + 2F(\beta^*)\right)}{TP}$$

*Where we are*

*OPT*

*# iters* — *TP* — *# procs*

*linear speed ups* *up to P processors*

Generalizes bounds for Shooting (Shalev-Shwartz & Tewari, 2009)

---

# Convergence Analysis

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

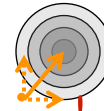<u>Theorem: Shotgun Convergence</u>

Assume $P < d/\rho + 1$

where $\rho$ = spectral radius of $\mathbf{X^T X}$

$$E\left[F(\beta^{(T)})\right] - F(\beta^*)$$
$$\leq \frac{d\left(\frac{1}{2} \| \beta^* \|_2^2 + F(\beta^{(0)})\right)}{TP}$$

Nice case:
Uncorrelated features

$\rho = \underline{1} \Rightarrow P_{max} = \underline{d}$

Bad case:
Correlated features

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

# Empirical Evaluation

Mug32_singlepixcam

Iterations to convergence vs P (# *simulated* parallel updates)

$P_{max}=158$

$d = 1024$

$\rho = 6.4967$

Ball64_singlepixcam

Iterations to convergence vs P (# *simulated* parallel updates)

$P_{max}=3$

$d = 4096$

$\rho = 2047.8$

©Carlos Guestrin 2013

17

---

# Stepping Back…

- Stochastic coordinate ascent  SCD
  - Optimization: Pick a coordinate j; find min $\beta_j$
  - Parallel SCD: Pick P coordinates
  - Issue: coordinates may interfere P coordinates
  - Solution: bound possible interference based on $\rho$
- Natural counterpart: SGD
  - Optimization: Pick a data points $\beta \in \beta - \eta \, \nabla F(x^i, \beta)$
  - Parallel: Pick P data points & indep update $\beta$
  - Issue: can interfere in all coordinates
  - Solution: bound interference

©Carlos Guestrin 2013

18

---

9

# Parallel SGD with No Locks

- Each processor in parallel:
  - □ Pick data point i at random
  - □ For j = 1…d:

$$\beta_j \leftarrow \beta_j - \eta \left( \nabla F(x^i, \beta) \right);$$

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

Other interferences

19

# Addressing Interference in Parallel SGD

- Key issues:
  - □ Old gradients

$DF(x^{i2}, \beta^{(7)})$ — $\Delta_\beta^{(1)}$

Processor ① $\beta^{(7)}$

$\beta^{(8)}$

② $\Delta_\beta^{(12)}$

  - □ Processors overwrite each other's work

- Nonetheless:
  - □ Can achieve convergence and some parallel speedups
  - □ Proof uses weak interactions, but through sparsity of data points
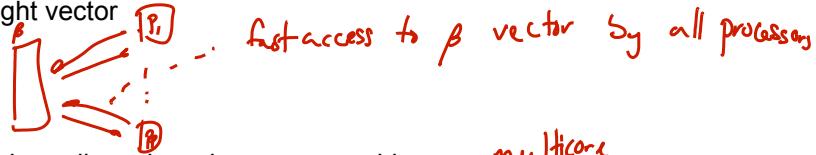
sparsity X
is key to analysis

20

10

# Problem with Parallel SCD and SGD

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

  *fast access to $\beta$ vector by all processors*

- Works well on shared memory machines   *multicore*

- Very difficult to implement efficiently in distributed memory, *Cloud*

  *$\beta := \leftarrow$ slow connection*  $P_1$ ... $P_2$

- Open problem: Good parallel SGD and SCD for distributed setting…
  - Let's look at a trivial approach

# Simplest Distributed Optimization Algorithm Ever Made

- Given $N$ data points & $P$ machines
- Stochastic optimization problem: $\min_\beta F(\beta) = \frac{1}{N} \sum_{i=1}^{N} F(x^i; \beta)$
- Distribute data: *P machines*

  *randomly assign data*  $P_1$ ... $P_P$

  *solves a problem $D_K$*  $|D_K| = \frac{N}{P} = n$

- Solve problems independently

  *machine $k$ independently estimates* $\beta^{(k)} = \min_\beta \frac{1}{n} \sum_{x^i \in D_K} F(x^i; \beta)$

- Merge solutions

  $$\tilde{\beta} = \frac{1}{P} \sum_k \beta^{(K)}$$

- Why should this work at all????

# For Convex Functions…
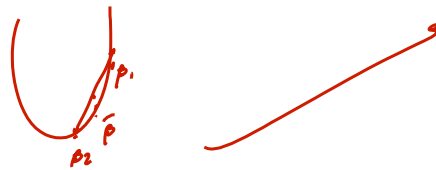
- Convexity:

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

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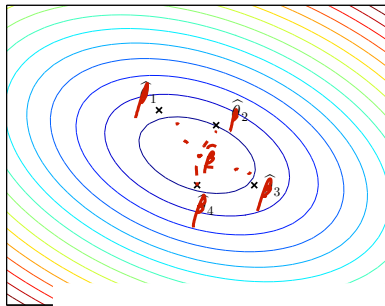
$$\max(F(\beta_1), F(\beta_2)) \geq F(\bar{\beta})$$

23

# Hopefully…

- Convexity only guarantees:

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

- But, estimates from independent data!

24

## Analysis of Distribute-then-Average
[Zhang et al. '12]

- Under some conditions, including strong convexity, lots of smoothness, and more…

- If all data were in one machine, converge at rate: $\hat{\beta}_N \in \underset{\beta}{\text{argmin}} \frac{1}{N} \sum_{i=1}^{N} F(x^i, \beta)$

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

- With $p$ machines converge at a rate:

$n = \frac{N}{p}$

unavoidable

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

e.g.    1T data points , 1000 machines    $p \simeq N^{\frac{1}{4}}$    $\frac{1}{N}$

=) plug in $\frac{1}{n^2}$ → negligible when compared to $N$

great at parallelism

25

---

# 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:

    Sparsity patterns in machine $i$    can be very different
    than those in machine $j$
    → average $\beta$ ⇒ 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

26

13

# What you need to know

- One way to solve LASSO problem
- Stochastic Coordinate Descent (SCD)
- Minimizing a coordinate in LASSO
- 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

27

14