**Case Study 3: fMRI Prediction**

# Fused LASSO
# LARS
# Parallel LASSO Solvers

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

Emily Fox

February 4th, 2014

1

---

# LASSO Regression

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

- New objective:

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

$$\underbrace{\phantom{\sum_{i=1}^{N} \left( y^i - (\beta_0 + \beta^T x^i) \right)^2}}_{RSS(\beta)}$$

$$\Updownarrow$$

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

2

1

# Geometric Intuition for Sparsity

Overall obj: $F(\beta) = RSS(\beta) + \lambda \|\beta\|$ ← 1 or 2 norm "lasso" "ridge"

solns typically hit corners

inc. RSS

$\hat{\beta}_\lambda^{lasso}$

$= \min F(\beta)$ for some $\lambda$

$\hat{\beta}_1 = 0$
$\hat{\beta}_2 = sthg$

level curves for $\|\beta\|_1 \leq B$

$\min \|\beta\|_1$

min RSS($\beta$)

$\hat{\beta}^{ML}$

$\beta_2$   $\beta_1$   Lasso

min RSS($\beta$)

$\hat{\beta}_\lambda^{ridge} = \min F(\beta)$ for some $\lambda$

$\hat{\beta}_1 = sthg$
$\hat{\beta}_2 = sthg$

$\beta_2$   $\beta_1$   Ridge Regression

# Soft Threshholding

For all other $\beta_k$ fixed, what must $\beta_j$ satisfy?
$c_j \propto corr(x_j, r_{-j})$

all examples of feature $j$

residual from model w/o using $j^{th}$ covariate

$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases}$$

$$= sign\left(\frac{c_j}{a_j}\right)\left(\frac{|c_j|}{a_j} - \frac{\lambda}{a_j}\right)_+$$

If $X^T X = I$

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

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

"is corr. strong enough?"

LS soln "ML"

ridge

lasso

From Kevin Murphy textbook

$\beta_j$   $c_j$   $-\lambda$   $\lambda$

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

## LASSO Coefficient Path

*Again, for each $\lambda$ we have a diff. soln*



From
Kevin Murphy
textbook

$\|\beta\|_1 \leq B$

*a few critical points where coeff. enter model*

*CV soln*

*B*

*inc. $\lambda$*

5

---

# Sparsistency

*typical*

- Typical Statistical Consistency Analysis:
  - Holding model size (*p*) fixed, as number of samples (*N*) goes to infinity, estimated parameter goes to true parameter

    *est param.* $\hat{\theta} \rightarrow \theta^{*}$ *true param* ?

- Here we want to examine *p >> N* domains
- Let both model size *p* and sample size *N* go to infinity!
  - Hard case: *N = k* log *p*

6

3

# Sparsistency

- Rescale LASSO objective by *N*:

$$\min_{\beta} \frac{1}{N} RSS(\beta) + \lambda_N \sum_j |\beta_j|$$

- Theorem (Wainwright 2008, Zhao and Yu 2006, …):
  - □ Under some constraints on the design matrix *X*, if we solve the LASSO regression using

$$\lambda_N > \frac{2}{\gamma} \sqrt{\frac{2\sigma^2 \log p}{N}}$$

    Then for some $c_1 > 0$, the following holds with at least probability

$$1 - 4\exp\left(-c_1 N \lambda_N^2\right) \longrightarrow 1$$

  - The LASSO problem has a unique solution with support contained within the true support
  - *(stronger:)* If $\min_{j \in S(\beta^*)} |\beta_j^*| > c_2 \lambda_N$ for some $c_2 > 0$, then $S(\hat{\beta}) = S(\beta^*)$
    
    ↑ *coeff large enough relative to penalty*

©Emily Fox 2014    7

---

# Case Study 3: fMRI Prediction

# Fused LASSO

Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
February 4th, 2014

©Emily Fox 2014    8

# fMRI Prediction Subtask

*If we use LASSO: penalizing individual voxels*

- **Goal:** Predict semantic features from fMRI image
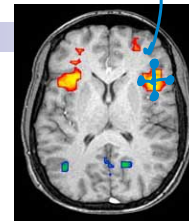
*Learning $S$: images $\rightarrow$ semantic features*



20,000 voxels

**Features of word**

$y^i$

$x^i$

*feature $j$*

Simplification: Consider predicting each $\hat{y}_j^i$ separately

$$x^i = \begin{bmatrix} \hat{x}_1^i \\ \hat{x}_2^i \\ \vdots \\ x_{20,000}^i \end{bmatrix} \xrightarrow{S} \begin{bmatrix} \hat{y}_1^i \\ \hat{y}_2^i \\ \vdots \\ \hat{y}_d^i \end{bmatrix}$$

$x^i \in \mathbb{R}^{20,000}$

$P$

$d = \#$ of semantic features

$y^i \in \mathbb{R}^d$

9

---

# Fused LASSO

*2d lattice*



- Might want coefficients of neighboring voxels to be similar

*discover important regions*

- How to modify LASSO penalty to account for this?

- Graph-guided fused LASSO
  - Assume a 2d lattice graph connecting neighboring pixels in the fMRI image
  - Penalty:

$$\|y - X\beta\|_2^2 + \lambda_1 \sum_j |\beta_j| + \lambda_2 \sum_{(s,t)\in E} |\beta_s - \beta_t|$$

$RSS(\beta)$

pair of voxels

edge set

penalizing diff. bt these weights

10

5

# Generalized LASSO

- Assume a structured linear regression model:

$$\|y - X\beta\|_2^2 + \lambda \|D\beta\|_1$$

$$D \in \mathbb{R}^{m \times p}$$

- If *D* is invertible, then get a new LASSO problem if we substitute

$$\beta = D^{-1}\beta^{new} \rightarrow \|y - XD^{-1}\beta^{new}\|_2^2 + \lambda \cdot \|\beta^{new}\|_1$$

$\underset{\tilde{X}}{\underbrace{\phantom{XD^{-1}}}}$ new design matrix

- Otherwise, not equivalent

  − solve for $\hat{\beta}^{new}$
  − set $\hat{\beta} = D^{-1}\hat{\beta}^{new}$

- For solution path, see
  Ryan Tibshirani and Jonathan Taylor, "The Solution Path of the
  Generalized Lasso." Annals of Statistics, 2011.

---

# Generalized LASSO

signal approximation scenario

$X = I$

each $y^i$ has a unique $x^i$

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 1 & 0 & 0 & \dots \\ 0 & -1 & 1 & 0 & \dots \\ 0 & 0 & -1 & 1 & \dots \\ \vdots & & & & \end{bmatrix}$. This is the 1d fused lasso.

$$\lambda \sum_j |\beta_{j+1} - \beta_j|$$

encouraging piecewise const.

# Generalized LASSO

$$\hat{\beta}_\lambda = \operatorname*{argmin}_{\beta \in \mathbb{R}^n} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Suppose $D$ gives "adjacent" differences in $\beta$:

$$D_i = (0, 0, \ldots -1, \ldots, 1, \ldots 0),$$

$\leftarrow s \quad \nwarrow t$

$\lambda \sum_{(s,t) \in E} |\beta_t - \beta_s|$

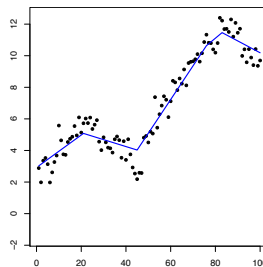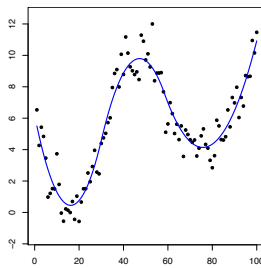where adjacency is defined according to a graph $\mathcal{G}$. For a 2d grid, this is the 2d fused lasso.

noisy image



true $\longrightarrow$

recovered image

13

---

# Generalized LASSO

$$\hat{\beta}_\lambda = \operatorname*{argmin}_{\beta \in \mathbb{R}^n} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 2 & -1 & 0 & \ldots \\ 0 & -1 & 2 & -1 & \ldots \\ 0 & 0 & -1 & 2 & \ldots \\ \vdots & & & & \end{bmatrix}$. This is linear trend filtering.



$\beta_3$

$\beta_2$

$\beta_1$

$\beta_3 - \beta_2 = \beta_2 - \beta_1$

$\Rightarrow 2\beta_2 - \beta_1 - \beta_3 = 0$

14

# Generalized LASSO

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 3 & -3 & 1 & \dots \\ 0 & -1 & 3 & -3 & \dots \\ 0 & 0 & -1 & 3 & \dots \\ \vdots & & & & \end{bmatrix}$. Get quadratic trend filtering.
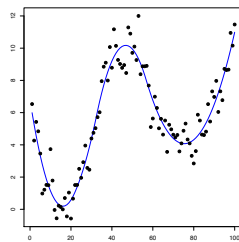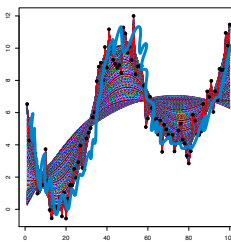
**15**

---

# Generalized LASSO

- Tracing out the fits as a function of the regularization parameter



$\hat{\beta}_\lambda$ for $\lambda = 25$        $\hat{\beta}_\lambda$ for $\lambda \in [0, \infty]$

**16**

# Acknowledgements

- Some material relating to the fused/generalized LASSO slides was provided by Ryan Tibshirani

17

---

## Case Study 3: fMRI Prediction

# LASSO Solvers – Part 1: LARS

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

Emily Fox
February 4th, 2014

18

# LASSO Algorithms

- Standard convex optimizer
- Now: Least angle regression (LAR)    *LARS = LAR + shrinkage*
  - □ Efron et al. 2004
  - □ Computes entire path of solutions
  - □ State-of-the-art until 2008
- Next up:
  - □ Pathwise coordinate descent ("shooting") – new
  - □ Parallel (approx.) methods

**19**

---

# LARS – Efron et al. 2004

- LAR is an efficient stepwise variable selection algorithm
  - □ "useful and less greedy version of traditional forward selection methods"

  *Efron*

- Can be modified to compute regularization path of LASSO
  - □ → LARS (Least angle regression and *shrinkage*)

- Increasing upper bound *B*, coefficients gradually "turn on"
  - □ Few critical values of *B* where support changes
  - □ Non-zero coefficients increase or decrease linearly between critical points
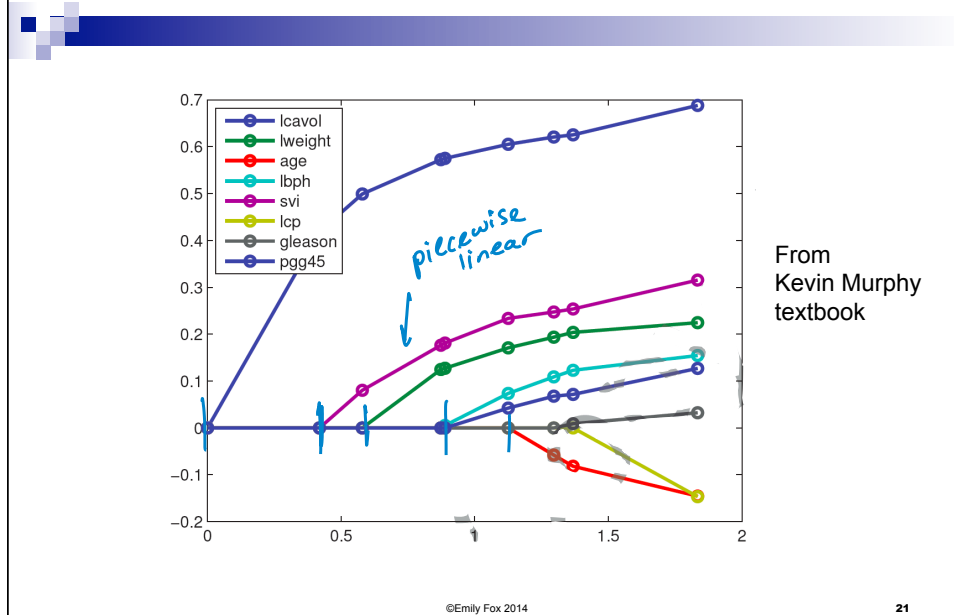  - ☆ □ Can solve for critical values analytically ←

- Complexity:

$$O\left(\min\left(N p^2,\, P N^2\right)\right)$$

*# obs*    *# of covariates*      *= cost of a single LS soln*

**20**

---

10

# LASSO Coefficient Path



From
Kevin Murphy
textbook

# LARS – Algorithm Overview

- Start with all coefficient estimates $\hat{\beta}_1 = \hat{\beta}_2 = \cdots \hat{\beta}_p = 0$

- Let $\mathcal{A}$ be the "active set" of covariates most correlated with the "current" residual $\leftarrow$ based on covariates already in model

- Initially, $\mathcal{A} = \{x_{j_1}\}$ for some covariate $x_{j_1}$

- Take the largest possible step in the direction of $x_{j_1}$ until another covariate $x_{j_2}$ enters $\mathcal{A}$

- Continue in the direction equiangular between $x_{j_1}$ and $x_{j_2}$ until a third covariate $x_{j_3}$ enters $\mathcal{A}$

- Continue in the direction equiangular between $x_{j_1}, x_{j_2}, x_{j_3}$ until a fourth covariate $x_{j_4}$ enters $\mathcal{A}$

- This procedure continues until all covariates are added at which point

# Comments

- LARS increases $\mathcal{A}$, but LASSO allows it to decrease

- Only involves a single index at a time

- If $p > N$, LASSO returns at most $N$ variables

- If group of variables are highly correlated, LASSO tends to choose one to include rather arbitrarily
  - Straightforward to observe from LARS algorithm....Sensitive to noise.

*beware of interpreting the variables included*

---

# More Comments

- In general, can't solve analytically for GLM (e.g., logistic reg.)
  - Gradually decrease λ and use efficiency of computing $\hat{\beta}(\lambda_k)$ from $\hat{\beta}(\lambda_{k-1})$ = warm-start strategy
  - See Friedman et al. 2010 for coordinate ascent + warm-starting strategy

- If $N > p$, but variables are correlated, ridge regression tends to have better predictive performance than LASSO (Zou & Hastie 2005)
  - Elastic net is hybrid between LASSO and ridge regression

$$\|y - X\beta\|_2^2 + \lambda_1 \sum |\beta_j| + \lambda_2 \|\beta\|_2^2$$

*(there still some issues... detail KM book)*

## Case Study 3: fMRI Prediction

LASSO Solvers – Part 2:
SCD for LASSO (Shooting)
Parallel SCD (Shotgun)
Parallel SGD
Averaging Solutions

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

Emily Fox

February 4th, 2014

**25**

# Scaling Up LASSO Solvers

- Another 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

**26**

# Coordinate Descent

- Given a function $F(\beta)$
  - Want to find minimum $\beta^* \in \min_\beta F(\beta) \leftarrow F(\beta_1, \cdots, \beta_p)$

- Often, hard to find minimum for all coordinates, but easy for one coordinate
  $1\text{-d optimization problem}$

- Coordinate descent:
  while not converged
    pick coord. $j$

  $\beta_j \leftarrow \min_b F(\beta_1, \beta_2 \cdots \beta_{j-1}, b, \beta_{j+1}, \cdots, \beta_p)$
                  varying $j^{th}$ coord. only

- How do we pick a coordinate?
  Round robin, random, smartly....

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

27

---

# Soft Threshholding

For all other coord. fixed, soln for $j^{th}$ coord
$c_j \propto \text{corr}(\underline{x}_j, \underline{r}_{-j})$

$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases}$$

$\min_b F(\cdots \beta_{j-1}, b, \beta_{j+1} \cdots)$

all examples of feature $j$

residual from model w/o using $j^{th}$ covariate

$= \text{sign}\left(\frac{c_j}{a_j}\right)\left(\frac{|c_j|}{a_j} - \frac{\lambda}{a_j}\right)_+ = \text{sign}(c_j)\left(\frac{|c_j| - \lambda}{a_j}\right)$

$\text{if } c_j \text{ corr strong enough?"}$

If $X^T X = I$

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

$\hat{\beta}_j^{ridge} = \frac{\hat{\beta}_j^{LS}}{1+\lambda}$

LS soln $= \hat{\beta}^{LS}$
ridge
lasso

$\beta_j$    $c_j$

From Kevin Murphy textbook

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

28

14

## 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:  *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)

Can be done more efficiently.   Proof: HW!

---

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

$e_j : \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \leftarrow j^{th}$ coord.

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

- For (coordinate-wise) strongly convex functions:
$$F(\beta + \Delta\beta) \leq F(\beta) + \partial\beta_j (\nabla F(\beta))_j + \frac{\gamma (\partial\beta_j)^2}{2}$$

$\uparrow \Delta\beta = \partial\beta_j \cdot e_j$

Lasso:  $\gamma = 1$

Log. reg.  $\gamma = \frac{1}{4}$

dim    how hard    where we start from

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

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

$\leftarrow$ gets linearly better w/ iters

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

- Natural question: How does SCD & SGD convergence rates differ?

See paper:   SCD $\rightarrow$ faster w/ larger $p$  $\leftarrow$ no params to tune
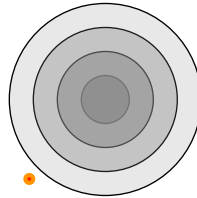SGD $\rightarrow$ faster w/ larger N  $\leftarrow$ needs $\eta$

# 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?
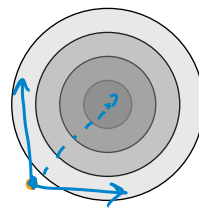  - do a round robin iter to msr convergence

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

Lasso: $\min_{\beta} F(\beta)$ where $F(\beta) = \| X\beta - \mathbf{y} \|_2^2 + \lambda \| \beta \|_1$
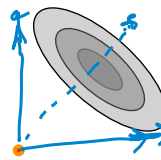
Shotgun (Parallel SCD)
While not converged,
- On each of $P$ processors,
  - Choose random coordinate j,
  - Update $\beta_j$ (same as for Shooting)

$\beta$ processors
$\beta_2$ ... 1
$\beta_7$ ... P
act as if they were the only children

yes!
features are uncorrelated

no !!
features are highly corr.

# Is SCD inherently sequential?

Lasso: $\min_\beta F(\beta)$ where $F(\beta) = \parallel X\beta - \mathbf{y} \parallel_2^2 + \lambda \parallel \beta \parallel_1$

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

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

33

---

# Is SCD inherently sequential?

Lasso: $\min_\beta F(\beta)$ where $F(\beta) = \parallel X\beta - \mathbf{y} \parallel_2^2 + \lambda \parallel \beta \parallel_1$

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 pos. or neg.

"interference" or "bias" from parallelism

34

17

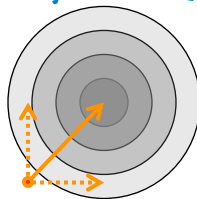# Is SCD inherently sequential?
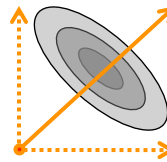
Theorem: If X is normalized s.t. diag($X^TX$)=1,

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

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

*key term* ← msrs magnitude of interference

$(X^TX)_{jk} = 0$
corr. bt. $X_j \sim X_k$

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

Nice case: Uncorrelated features

Bad case: Correlated features

©Emily Fox 2014    35

---

# Shotgun: Convergence Analysis

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

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

↑dim   spectral radius of $X^TX$

$$E\left[F(\beta^{(T)})\right] - F(\beta^*) \leq \frac{\overset{dim}{p}\left(\|\beta^*\|_2^2 + 2F(\beta^{(0)})\right)}{T \cdot P}$$

where we are

opt.

← # of processors

↑ # iters

linear speed up up to P proc.

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

©Emily Fox 2014    36

18

# Convergence Analysis

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

Theorem: Shotgun Convergence

Assume $P < p/\rho + 1$

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

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

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

Nice case:
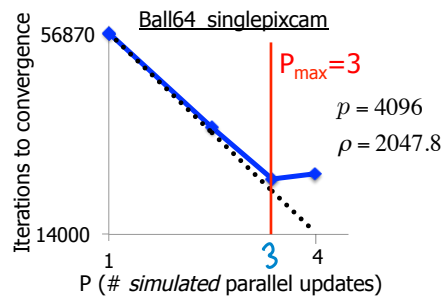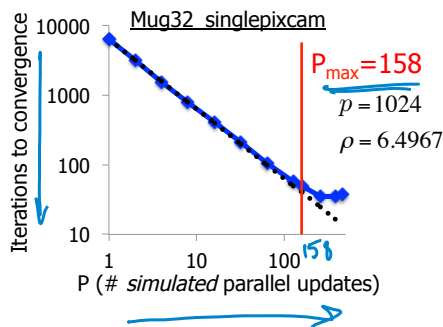Uncorrelated features

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

Bad case:
Correlated features

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

37

---

# Empirical Evaluation

2 classical compressed sensing datasets

Mug32_singlepixcam

Iterations to convergence

$P_{max}$=158

$p = 1024$

$\rho = 6.4967$

P (# *simulated* parallel updates)

158

Ball64_singlepixcam

Iterations to convergence

56870

$P_{max}$=3

$p = 4096$

$\rho = 2047.8$

14000

P (# *simulated* parallel updates)

3

38

19

# Stepping Back…

- Stochastic coordinate ascent  *SCD*
  - 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

39

---

# 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…*p*:

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

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

  other interferences

40

# What you need to know

- Sparsistency
- Fused LASSO
- LASSO Solvers
  - LARS
  - A simple SCD for LASSO (Shooting)
    - Your HW, a more efficient implementation! ☺
    - Analysis of SCD
  - Parallel SCD (Shotgun)

41