

Case Study 3: fMRI Prediction

LASSO Review, Fused LASSO, Parallel LASSO Solvers

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

Emily Fox

April 28th, 2015

©Emily Fox 2015

1

LASSO Regression

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

- New objective:

$$\min_{\beta} \underbrace{\sum_{i=1}^n (y^i - (\beta_0 + \beta^T x^i))^2}_{RSS(\beta)} + \lambda \|\beta\|_1$$

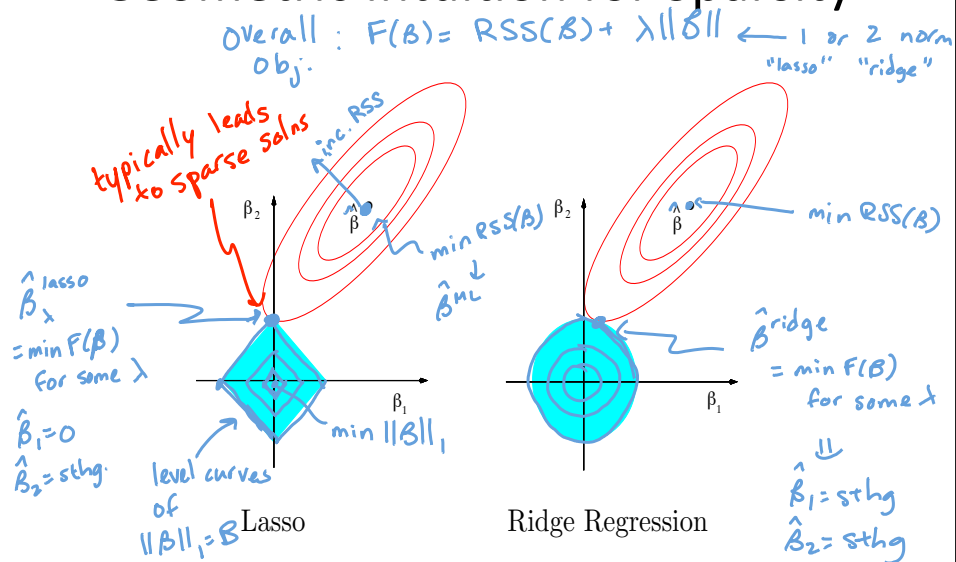
← L_1 penalty

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

©Emily Fox 2015

2

Geometric Intuition for Sparsity



©Emily Fox 2015

3

If we hold β_j fixed, then this must hold:

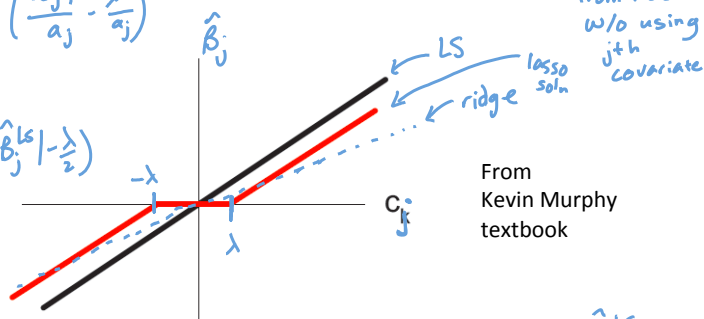
$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases}$$

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

If $X^T X = I$

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

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

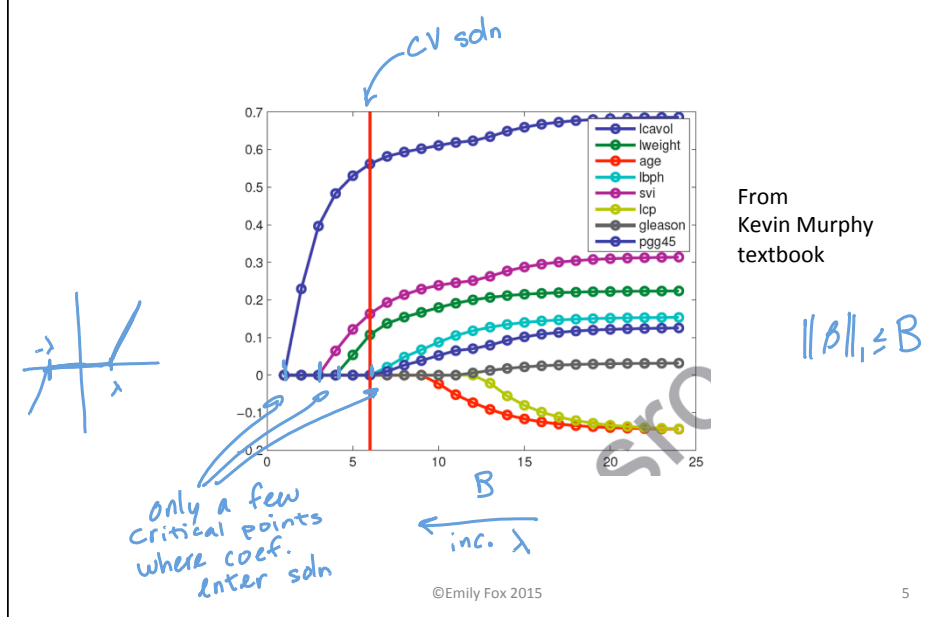


In LASSO, all coeff. are shrunk relative to $\hat{\beta}^{\text{ls}}$

©Emily Fox 2015

4

LASSO Coefficient Path



LASSO Example

CV soln

Term	Least Squares	Ridge	Lasso
Intercept	2.465	2.452	2.468
lcavol	0.680	0.420	0.533
lweight	0.263	0.238	0.169
age	-0.141	<u>-0.046</u>	
lbph	0.210	0.162	0.002
svi	0.305	0.227	0.094
lcp	-0.288	<u>0.000</u>	
gleason	-0.021	<u>0.040</u>	
pgg45	0.267	<u>0.133</u>	

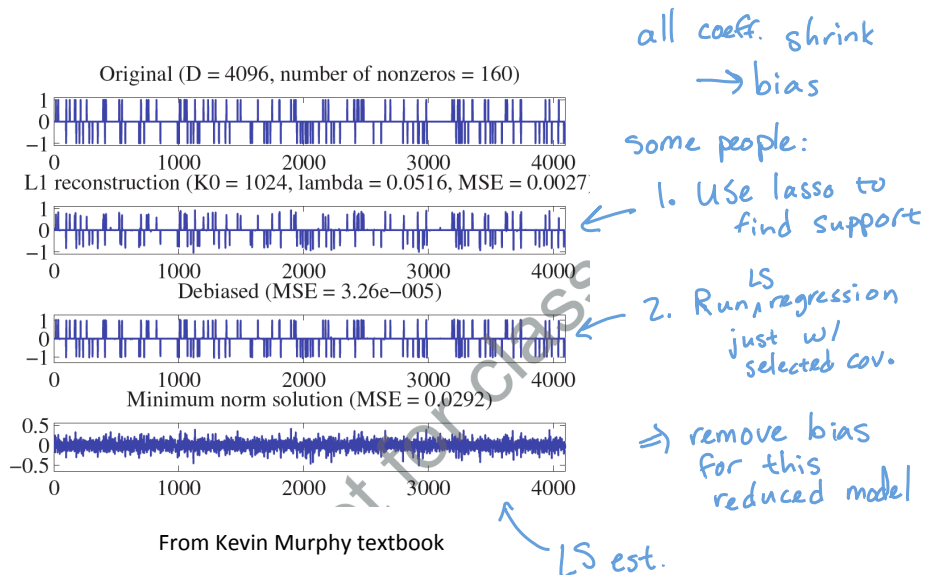
not in the model

shrunk, but non-zero

©Emily Fox 2015

6

Debiasing



©Emily Fox 2015

7

Sparsistency

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

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

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

N grows slowly relative to p

©Emily Fox 2015

8

Sparsistency

- Rescale LASSO objective by N :

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

- 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

$$1 - 4 \exp(-c_1 N \lambda_N^2) \rightarrow 1$$

- The LASSO problem has a unique solution with support contained within the true support

$$S(\hat{\beta}) \subseteq S(\beta^*)$$

- If $\min_{j \in S(\beta^*)} |\beta_j^*| > c_2 \lambda_N$ for some $c_2 > 0$, then $S(\hat{\beta}) = S(\beta^*)$

\sim coeff large enough relative to penalty

©Emily Fox 2015

9

Case Study 3: fMRI Prediction

Fused LASSO

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

Emily Fox

April 28th, 2015

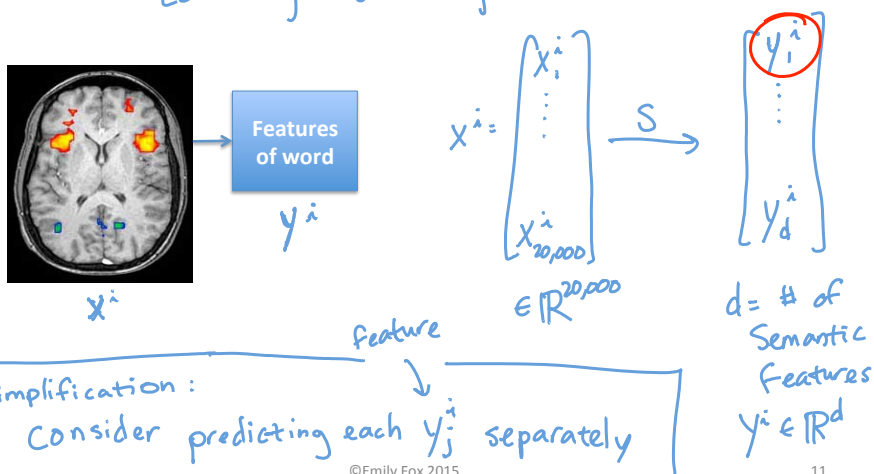
©Emily Fox 2015

10

fMRI Prediction Subtask

- **Goal:** Predict semantic features from fMRI image

Learning $S: \text{images} \rightarrow \text{semantic features}$

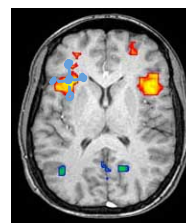


Fused LASSO

- Might want coefficients of neighboring voxels to be similar

discover regions of importance

- 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:

$$\|y - XB\|_2^2 + \lambda_1 \sum_j |B_j| + \lambda_2 \sum_{(s,t) \in E} |B_s - B_t|$$

penalizing coeffs. having different weights

$(s,t) \in E$ has edge in graph

©Emily Fox 2015

12

Generalized LASSO

- Assume a structured linear regression model:

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

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

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

$$\beta = D^{-1}\beta^{\text{new}} \rightarrow \|y - \underbrace{XD^{-1}}_{\text{new design matrix}}\beta^{\text{new}}\|_2^2 + \lambda \|\beta^{\text{new}}\|_1$$

- Otherwise, not equivalent

- For solution path, see

Ryan Tibshirani and Jonathan Taylor, "The Solution Path of the Generalized Lasso." Annals of Statistics, 2011.

©Emily Fox 2015

13

Generalized LASSO

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

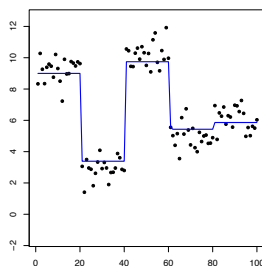
"signal approximation"
scenario $X=I$

- associate a
unique x_i
w/ each y_i

$$\text{Let } D = \begin{bmatrix} -1 & 1 & 0 & 0 & \dots \\ 0 & -1 & 1 & 0 & \dots \\ 0 & 0 & -1 & 1 & \dots \\ \vdots & & & & \ddots \end{bmatrix}$$

This is the **1d fused lasso**.

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



encourages
piecewise
constant

©Emily Fox 2015

14

Generalized LASSO

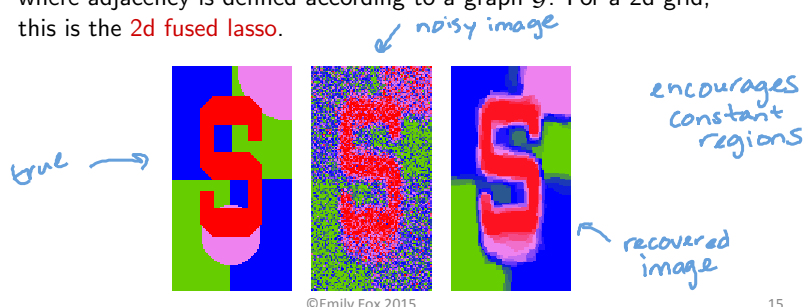
$$\hat{\beta}_\lambda = \operatorname{argmin}_{\beta \in \mathbb{R}^n} \frac{1}{2} \|y - \beta\|_2^2 + \lambda \|D\beta\|_1$$

Suppose D gives “adjacent” differences in β :

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

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

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

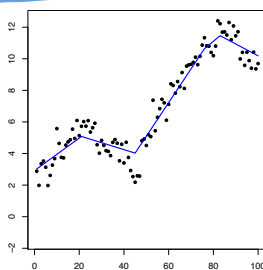


15

Generalized LASSO

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

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

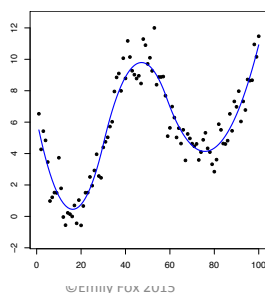


16

Generalized LASSO

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

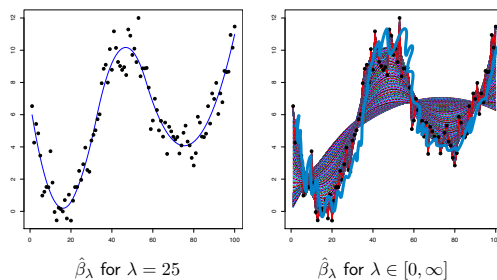
Let $D = \begin{bmatrix} -1 & 3 & -3 & 1 & \dots \\ 0 & -1 & 3 & -3 & \dots \\ 0 & 0 & -1 & 3 & \dots \\ \vdots & & & & \end{bmatrix}$. Get quadratic trend filtering.



17

Generalized LASSO

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



©Emily Fox 2015

18

Acknowledgements

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

©Emily Fox 2015

19

Case Study 3: fMRI Prediction

LASSO Solvers—Part 1: LARS

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

Emily Fox

April 28th, 2015

©Emily Fox 2015

20

LASSO Algorithms

- So far: Standard convex optimizer
- Now: Least angle regression (LAR)
 - 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

LARS ← shrinkage
 } now
 } later

©Emily Fox 2015

21

LARS – Efron et al. 2004

- LAR is an efficient stepwise variable selection algorithm
 - “useful and less greedy version of traditional forward selection methods”
- 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:

$$O(\min(Np^2, pN^2))$$

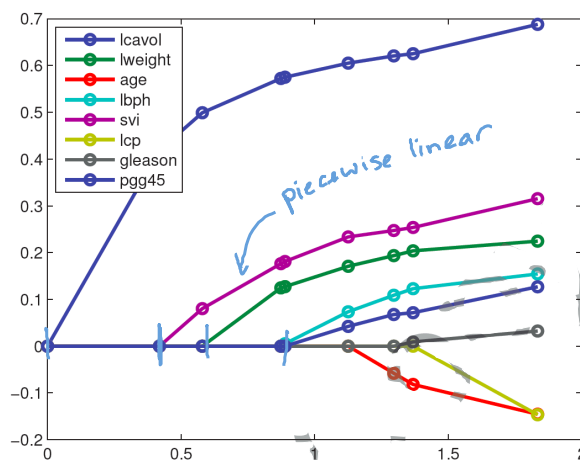
\nwarrow # of obs \nearrow # of covariates = cost of a single LS soln

key to providing full reg. path/coeff.

©Emily Fox 2015

22

LASSO Coefficient Path



From
Kevin Murphy
textbook

©Emily Fox 2015

23

LARS – Algorithm Overview

- Start with all coefficient estimates $\hat{\beta}_1 = \hat{\beta}_2 = \dots = \hat{\beta}_p = 0$
- Let \mathcal{A} be the “active set” of covariates most correlated with the “current” residual
← based on covariates already in the 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

©Emily Fox 2015

24

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

©Emily Fox 2015

25

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-start 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 are still some issues ...
see KM book)

©Emily Fox 2015

26

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

April 28th, 2015

©Emily Fox 2015

27

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
- ADMM

©Emily Fox 2015

28

Coordinate Descent

- Given a function $F(\beta)$
 - Want to find minimum $\beta^* \leftarrow \arg\min_{\beta} F(\beta)$ $\leftarrow F(\beta_1, \dots, \beta_p)$
- Often, hard to find minimum for all coordinates, but easy for one coordinate
1D optimization problem
- Coordinate descent:
 - while not converged
 - pick coordinate j \leftarrow varying j th coord.
 - $\beta_j \leftarrow \arg\min_b F(\beta_1, \beta_2, \dots, \beta_{j-1}, b, \beta_{j+1}, \dots, \beta_p)$
- How do we pick a coordinate?
 - Round robin, random, smartly, ...
- When does this converge to optimum?
 - e.g. strongly convex (separability)

©Emily Fox 2015

29

Soft Thresholding

Fixing all other coord., soln for j th coord:

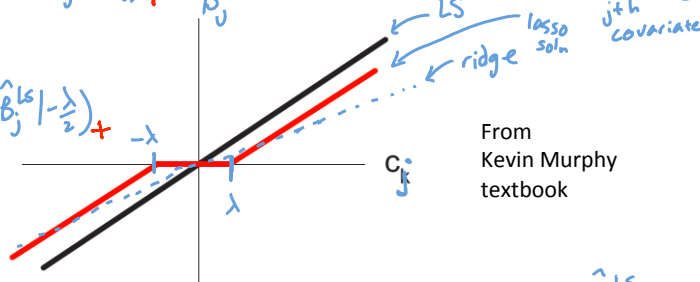
$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases}$$

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

If $X^T X = I$

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

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



From Kevin Murphy textbook

In LASSO, all coeff. are shrunk relative to $\hat{\beta}^{\text{ls}}$

©Emily Fox 2015

30

Stochastic Coordinate Descent for LASSO (aka Shooting Algorithm)

- Repeat until convergence

– Pick a coordinate j at random

- Set:

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

- Where:

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

cache \nearrow *Cost per iteration* \nearrow $O(N)$ \nwarrow *omitting j's cov* $[\beta_1, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_p]$

Can be done more efficiently, Proof: Your HW!

©Emily Fox 2015

31

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

$$e_j = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow j^{\text{th}} \text{ dim}$$

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

- For (coordinate-wise) strongly convex functions:

$$F(B + \Delta B) \leq F(B) + \partial \beta_j (\nabla F(B))_j + \frac{\gamma (\partial \beta_j)^2}{2}$$

$$\Delta B = \partial \beta_j \cdot e_j$$

- Theorem:

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

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

dim \nwarrow *how hard* \nwarrow *where we started from*

Lasso:
 $\gamma = 1$

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

~~gets~~ gets linearly better with iterations

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

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

see paper: SCD \rightarrow faster w/ large P \leftarrow no params to tune
SGD \rightarrow faster w/ large N \leftarrow need η

©Emily Fox 2015

32

Shooting: Sequential SCD

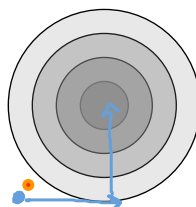
Lasso: $\min_{\beta} F(\beta)$ where $F(\beta) = \|X\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1$

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

While not converged,

- Choose random coordinate j ,
- Update β_j (closed-form minimization)

$F(\beta)$ contour



How do we measure?

- annoying
- over a time window? has anything changed?
* - do a round robin iter to measure convergence

©Emily Fox 2015

33

Shotgun: Parallel SCD [Bradley et al '11]

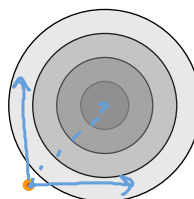
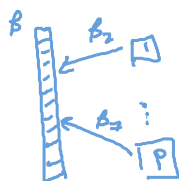
Lasso: $\min_{\beta} F(\beta)$ where $F(\beta) = \|X\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_1$

Shotgun (Parallel SCD)

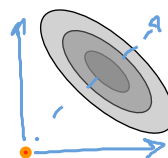
While not converged,

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

independently



yes!
features
are
uncorrelated



no!!
Feature
are highly
correlated

©Emily Fox 2015

34

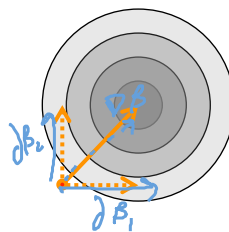
Is SCD inherently sequential?

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 update if
features are correlated.
Can we quantify this?

©Emily Fox 2015

35

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. $\text{diag}(X^T X) = 1$,

$$F(\beta + \Delta\beta) - F(\beta) \leq - \underbrace{\sum_{i_j \in \mathcal{P}} (\delta\beta_{i_j})^2}_{\text{"positive" progress}} + \sum_{\substack{i_j, i_k \in \mathcal{P}, \\ j \neq k}} \underbrace{(X^T X)_{i_j, i_k}}_{\text{could be pos. or neg.}} \delta\beta_{i_j} \delta\beta_{i_k}$$

"interference" or "bias" from parallelism

©Emily Fox 2015

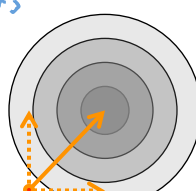
36

Is SCD inherently sequential?

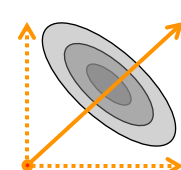
Theorem: If X is normalized s.t. $\text{diag}(X^T X) = 1$,

$$F(\beta + \Delta\beta) - F(\beta) \leq - \sum_{i_j \in \mathcal{P}} (\delta\beta_{i_j})^2 + \sum_{\substack{i_j, i_k \in \mathcal{P}, \\ j \neq k}} \underbrace{(X^T X)_{i_j, i_k}}_{\text{key term} \leftarrow \text{snsrs magnitude of interference}} \delta\beta_{i_j} \delta\beta_{i_k}$$

$(X^T X)_{jk} = 0$ \leftarrow coll bet. x_j & x_k
 $(X^T X)_{jk} \neq 0$ \leftarrow "interference"



Nice case:
Uncorrelated
features



Bad case:
Correlated
features

©Emily Fox 2015 37

Shotgun: Convergence Analysis

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

Assume # parallel updates $P < \overset{\text{dim}}{\rho} / \rho + 1$
 \leftarrow spectral radius $(X^T X)$

$$\underbrace{E[F(\beta^{(T)})]}_{\text{where we are}} - \underbrace{F(\beta^*)}_{\text{opt}} \leq \frac{\overset{\text{dim}}{P} (\|\beta^*\|_2^2 + 2F(\beta^{(0)}))}{\underset{\substack{\text{\# of iters} \\ \uparrow}}{T \cdot P} \leftarrow \text{\# of processors}}$$

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

©Emily Fox 2015

38

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

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

↑ linear speed up,
up to P processors

Nice case:
Uncorrelated
features



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

Bad case:
Correlated
features



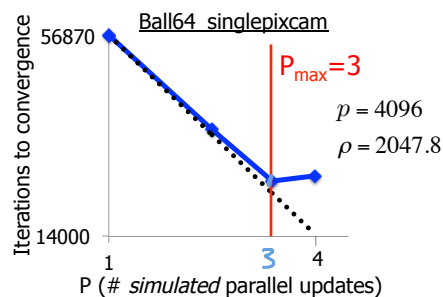
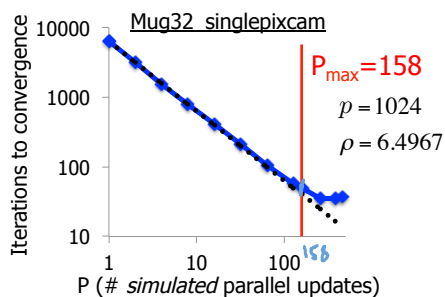
$$\rho = p \Rightarrow P_{\max} = 1 \text{ (at worst)}$$

©Emily Fox 2015

39

Empirical Evaluation

2 classical compressed
sensing datasets



©Emily Fox 2015

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)