

## Module 1: Nonparametric Preliminaries

### LASSO cont'd

STAT/BIOSTAT 527, University of Washington

Emily Fox

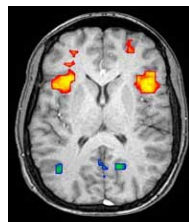
April 7<sup>th</sup>, 2014

©Emily Fox 2014

1

## fMRI Prediction Subtask

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



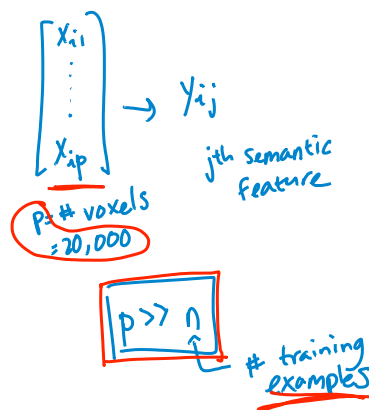
Features  
of word

$y_i$

$$\hat{\beta}^{LS} = (X^T X)^{-1} X^T y$$

Diagram showing the matrix dimensions:  $X$  is  $n \times p$  and  $y$  is  $n \times 1$ .

rank  
deficient



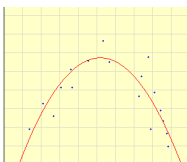
©Emily Fox 2014

2

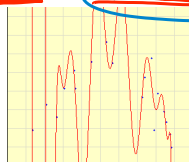
# Regularization in Linear Regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$



$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



even for  
 $n > p$ ,  
 $p$  large

- Regularized** or **penalized** regression aims to impose a “complexity” penalty by penalizing large weights
  - “Shrinkage” method

©Emily Fox 2014

3

# Ridge Regression

- Ameliorating issues with overfitting: *penalization of weights “regularization”*

- New objective: *LS obj.*

$$\min_{\beta} \sum_{i=1}^n (y_i - (\beta_0 + \beta^T x_i))^2 + \lambda \|\beta\|_2^2$$

don't penalize intercept

$\lambda =$  strength of penalty

$$\min_{\beta} \text{RSS}(\beta) \quad \text{s.t.} \quad \|\beta\|_2^2 \leq S$$

©Emily Fox 2014

4

# Variable Selection

- Ridge regression: Penalizes large weights

- What if we want to perform "feature selection"? *variable*

- E.g., Which regions of the brain are important for word prediction?
- Can't simply choose predictors with largest coefficients in ridge solution
- Computationally impossible to perform "all subsets" regression

*discrete*

*2<sup>P</sup> subsets of predictors ... can't do this*

- Stepwise procedures are sensitive to data perturbations and often include features with negligible improvement in fit

*greedy, w/ backtracking alg.*

- Try new penalty: Penalize non-zero weights

- Penalty:

$$L_1 \quad \|\beta\|_1 = \sum_j |\beta_j|$$

- Leads to sparse solutions
- Just like ridge regression, solution is indexed by a continuous param  $\lambda$

*not min this obj.  
- coeff. sensitive to what's in the model*

©Emily Fox 2014

5

# LASSO Regression

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

- New objective:

$$\min_{\beta} \underbrace{\sum_{i=1}^n (y_i - (\beta_0 + \beta^T x_i))^2}_{\text{RSS}(\beta)} + \underbrace{\lambda \|\beta\|_1}_{\text{L}_1 \text{ penalty}}$$

*LS obj.*  *$\sum |\beta_j|$*

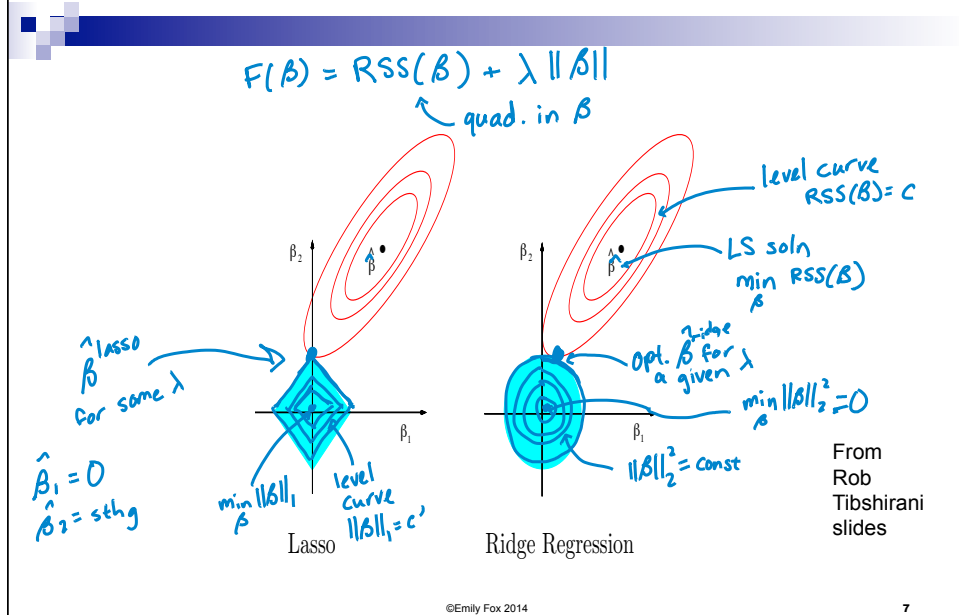


$$\min_{\beta} \text{RSS}(\beta) \quad \text{s.t.} \quad \|\beta\|_1 \leq B$$

©Emily Fox 2014

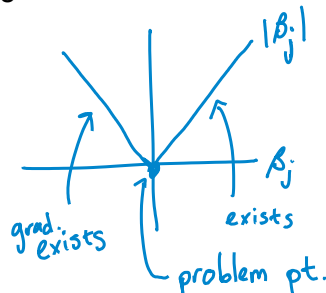
6

# Geometric Intuition for Sparsity



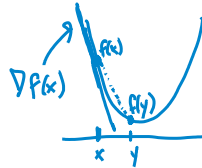
## Soft Thresholding

- To see why LASSO results in sparse solutions, look at conditions that must hold at optimum  
 look at  $\beta_j$  ... do this for all  $j \Rightarrow$  set of simultaneous equations
- $L_1$  penalty  $\|\beta\|_1$  is not differentiable whenever  $\beta_j = 0$   
 $\sum |\beta_j|$
- Look at subgradient...



# Subgradients of Convex Functions

- Gradients lower bound convex functions:

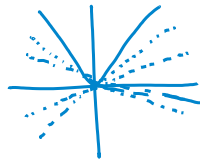


$$\frac{f(y) - f(x)}{y - x} \geq \nabla f(x)$$

$$\Rightarrow f(y) \geq f(x) + \nabla f(x)(y - x)$$

- Gradients are unique at  $x$  if function differentiable at  $x$
- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



For  $\beta_j$ :  
 $\forall \epsilon \in [-1, 1]$

$$\begin{aligned} &\forall \epsilon \in \partial f(x) \text{ subgrad.} \\ &\text{if} \\ &f(y) \geq f(x) + \epsilon(y - x) \end{aligned}$$

©Emily Fox 2014

9

# Soft Thresholding

Goal.  $\nabla_{\beta_j} (RSS(\beta) + \lambda \|\beta\|_1) = 0$

- Gradient of RSS term:

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

$$\uparrow 2 \sum_{i=1}^n (x_{ij})^2$$

$$2 \sum_{i=1}^n x_{ij} (y_i - \beta_j^T x_{i,-j})$$

all  $\beta$ 's except for  $\beta_j$       all cov's other than  $x_{ij}$

- Subgradient of full objective:

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

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

$c_j$  & corr ( $x_j, r_{-j}$ )  
msr of how relevant  $x_j$  is for pred  $y$  beyond what others can  
residuals from a model w/o  $j^{th}$  cov.

©Emily Fox 2014

10

# Soft Thresholding

- Set subgradient = 0:

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

If  $\beta_j < 0$

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

$$\Rightarrow \beta_j = \frac{c_j + \lambda}{a_j} < 0 \Rightarrow c_j < -\lambda \quad \text{If strong neg. corr., then } \hat{\beta}_j < 0$$

If  $\beta_j > 0$

$$a_j \beta_j - c_j + \lambda = 0 \Rightarrow \beta_j = \frac{c_j - \lambda}{a_j} > 0 \Rightarrow c_j > \lambda \quad \text{If strong pos. corr., then } \hat{\beta}_j > 0$$

If  $\beta_j = 0$   $-\lambda < c_j < \lambda$  if not strong corr., then  $\hat{\beta}_j = 0$

- The value of  $c_j = 2 \sum_{i=1}^N x_{ij}^T (y_i^T - \beta'_{-j} x_{i,-j}^T)$  constrains  $\beta_j$

©Emily Fox 2014

11

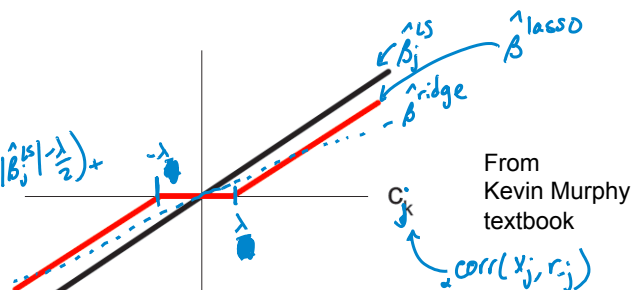
# Soft Thresholding

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

If  $X^T X = I$

$$\hat{\beta}_{\text{ridge}} = \frac{\hat{\beta}_{\text{LS}}}{1 + \lambda}$$

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



From Kevin Murphy textbook

In lasso, all coeff.  $\hat{\beta}_{\text{lasso}}$  are shrunk relative to  $\hat{\beta}_{\text{LS}}$

©Emily Fox 2014

12

# Coordinate Descent

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

©Emily Fox 2014

13

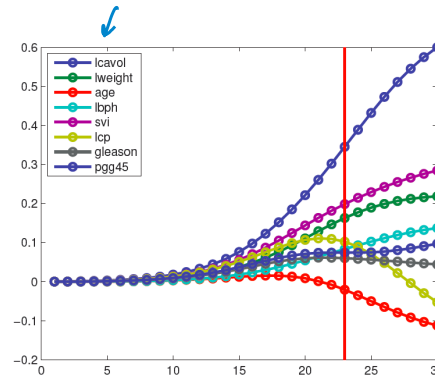
## Stochastic Coordinate Descent for LASSO (aka Shooting Algorithm)

- 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:  $c_j = 2 \sum_{i=1}^N x_{ij} (y_i - \beta'_{-j} x_{-j})$
- For convergence rates, see Shalev-Shwartz and Tewari 2009
- Other common technique = LARS
  - Least angle regression and shrinkage, Efron et al. 2004

©Emily Fox 2014

14

## Recall: *Ridge Coefficient Path*



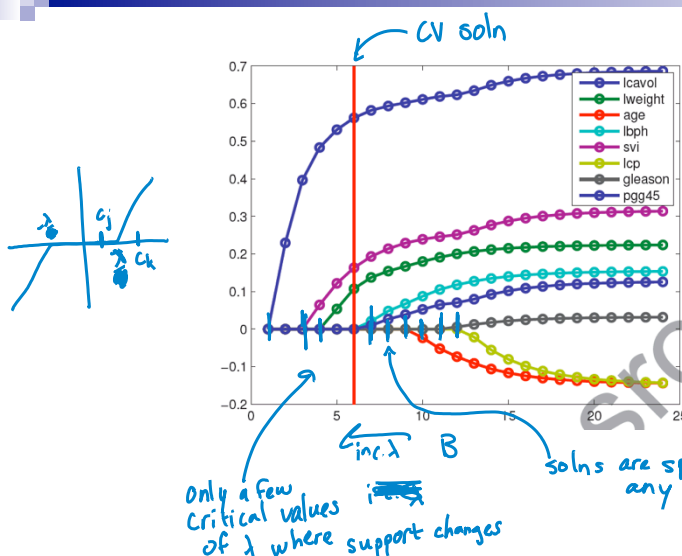
From  
Kevin Murphy  
textbook

- Typical approach: select  $\lambda$  using cross validation

©Emily Fox 2014

15

## Now: *LASSO Coefficient Path*



From  
Kevin Murphy  
textbook

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

©Emily Fox 2014

16



# LASSO Example

	Term	Least Squares	Ridge	Lasso
$\beta_0$	Intercept	2.465	2.452	2.468
$\beta_1$	lcavol	0.680	0.420	0.533
$\beta_2$	lweight	0.263	0.238	0.169
$\cdot$	age	-0.141	-0.046	
$\cdot$	lbph	0.210	0.162	0.002
$\cdot$	svi	0.305	0.227	0.094
	lcp	-0.288	0.000	
	gleason	-0.021	0.040	
$\beta_7$	pgg45	0.267	0.133	

From  
Rob  
Tibshirani  
slides

not in  
the model  
(sparse solns)

©Emily Fox 2014

17

# Sparsistency

## Typical Statistical Consistency Analysis:

- Holding model size ( $p$ ) fixed, as number of samples ( $n$ ) goes to infinity, estimated parameter goes to true parameter

$$\hat{\theta} \rightarrow \theta^* \text{ true param as } n \rightarrow \infty$$

- 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 2014

19

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

©Emily Fox 2014

20

# Comments

- In general, can't solve analytically for GLM (e.g., logistic reg.)

- Gradually decrease  $\lambda$  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

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

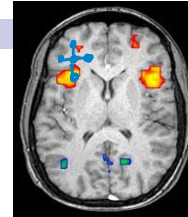
(still some issues, but other solns)

©Emily Fox 2014

21

# 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:

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

$\nwarrow$  *penalizing  $\beta_s \neq \beta_t$*        $\nearrow$  *in edge set*

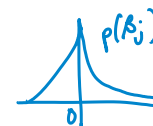


©Emily Fox 2014

22

# A Bayesian Formulation

- Consider a model with likelihood  
 $y_i | \beta \sim N(\beta_0 + x_i^T \beta, \sigma^2)$
- and prior  
 $\beta_j \sim \text{Lap}(\beta_j; \lambda)$



where  $\text{Lap}(\beta_j; \lambda) = \frac{\lambda}{2} e^{-\lambda |\beta_j|}$

- For large  $\lambda$   
*more peaked around 0*
- LASSO solution is equivalent to the **mode** of the posterior
- Note: posterior mode  $\neq$  posterior mean in this case  
*any given posterior sample is not sparse, but it will be penalized like in ridge.*
- There is no closed-form for the posterior. Rely on approx. methods.  
*spike+slab priors as alternatives*

$p(\beta_j = 0) = 0$   
*but posterior mode*

©Emily Fox 2014

23

# Reading

- Hastie, Tibshirani, Friedman: 3.4, 3.8.6

# What you should know

- LASSO objective
- Geometric intuition for differences between ridge and LASSO solns
- How LASSO performs soft thresholding
- Shooting algorithm
- Idea of sparsistency
- Ways in which other L1 and L1-Lp objectives can be encoded
  - Elastic net
  - Fused LASSO