# Module 1: Nonparametric Preliminaries

## LASSO cont'd

STAT/BIOSTAT 527, University of Washington

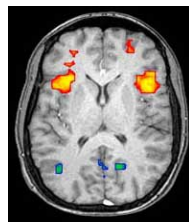Emily Fox

April 9th, 2013

---

# fMRI Prediction Subtask
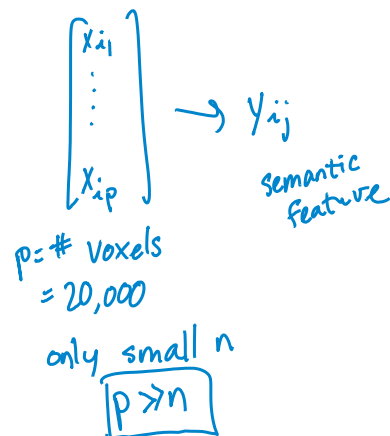
- **Goal:** Predict semantic features from fMRI image



Features of word

$y_i$

$x_i$

$\begin{bmatrix} x_{i_1} \\ \vdots \\ x_{i_p} \end{bmatrix} \rightarrow y_{ij}$

semantic feature

$p = \text{\# voxels} = 20,000$

only small $n$

$\boxed{p \gg n}$

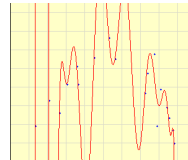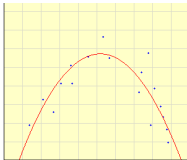$\hat{\beta}^{LS} = (X^T X)^{-1} X^T y$

$n$

rank deficient

# Regularization in Linear Regression

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

  $-2.2 + 3.1\ X - 0.30\ X^2$          $-1.1 + 4{,}700{,}910.7\ X - 8{,}585{,}638.4\ X^2 + \ldots$

  *even for*
  *$n > P$,*
  *$P$ large*

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

3

---

# Ridge Regression

- Ameliorating issues with overfitting:   *penalization of weights*
  *"regularization"*

- New objective:

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

  *LS.*
  *obj.*          *don't penalize the intercept*          *$\beta^T\beta$*

  *strength of the penalty*

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

4

2

# Variable Selection

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

- 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*

  $$2^P \text{ subsets of predictors} \dots \quad \text{can't do this}$$

  - Stepwise procedures are sensitive to data perturbations and often include features with negligible improvement in fit  $\leftarrow$ *greedy, ∄ backtracking alg.*

- Try new penalty: Penalize non-zero weights
  - Penalty:

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

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

5

# LASSO Regression

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

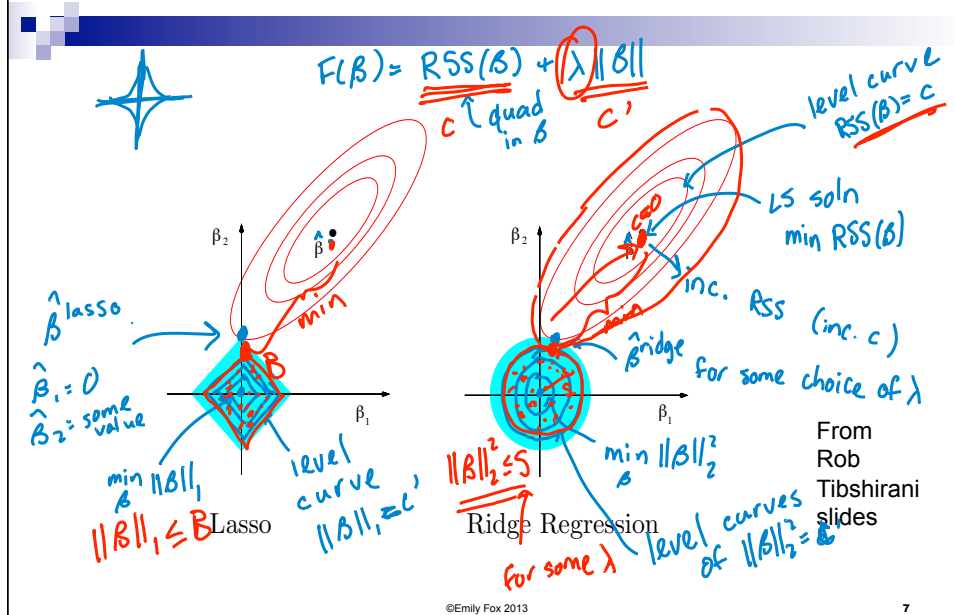- New objective:

$$\min_\beta \underbrace{\sum_{i=1}^n \left( y_i - (\beta_0 + \beta^T X_i) \right)^2}_{RSS(\beta)} + \boxed{\lambda \|\beta\|_1} \quad \bigstar$$

$$\Updownarrow$$

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

6

3

# Geometric Intuition for Sparsity

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

$c \uparrow$ quad in $\beta$   $c'$

level curve RSS($\beta$) = c

$\beta_2$

$\hat{\beta}$

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

$\hat{\beta}_1 = 0$
$\hat{\beta}_2 = $ some value

min

$\beta$

min $||\beta||_1$
$\beta$

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

level curve $||\beta||_1 = c'$

$\beta_1$

LS soln min RSS($\beta$)

inc. RSS (inc. c)

$\hat{\beta}^{ridge}$ for some choice of $\lambda$

$||\beta||_2^2 \leq S$

min $||\beta||_2^2$
$\beta$

Ridge Regression

for some $\lambda$   level curves of $||\beta||_2^2 = c'$

From Rob Tibshirani slides

©Emily Fox 2013    7

---

# Soft Threshholding

- To see why LASSO results in sparse solutions, look at conditions that must hold at optimum

  look at $\beta_j$ ... do this for all $j$ ⟹ set of simult. eqns

- $L_1$ penalty $||\beta||_1$ is not differentiable whenever $\beta_j = 0$

  $\sum |\beta_j|$

- Look at subgradient…

  $|\beta_j|$

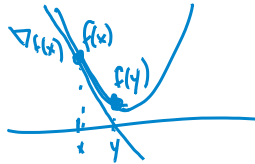  exists    problem pt    grad. exists    $\beta_j$

©Emily Fox 2013    8
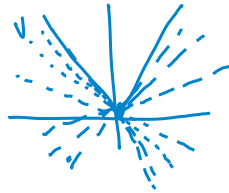
4

# Subgradients of Convex Functions

- Gradients lower bound convex functions:



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

- Gradients are <u>unique</u> at **x** if function differentiable at **x**

- Subgradients: Generalize gradients to non-differentiable points:
  - □ Any plane that lower bounds function:



For $|\beta_j|$:
$V \in [-1, 1]$

$$V \in \partial f(x) \text{ subgradient}$$
$$\text{if}$$
$$f(y) \geq f(x) + V(y-x)$$

9

---

# Soft Threshholding

Goal:
$$\nabla_{\beta_j} \left( \underbrace{RSS(\beta) + \lambda \|\beta\|_1}_{F(\beta)} \right) = 0$$

- Gradient of RSS term:

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

$$2 \sum_{n=1}^{n} x_j^i \left( y_i - \beta_{-j}^T x_{-j}^i \right) \quad \text{all cov other than } x_j$$
$$\text{all } \beta\text{'s except for } \beta_j$$

$$\uparrow \quad 2\sum_{n=1}^{n} (x_j^i)^2 \quad \text{Here: } x_{ij} = x_j^i$$

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

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

residuals from model w/o jth covariate

- Subgradient of full objective:

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

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

10

5

# Soft Threshholding

- 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 $\beta_j < 0$

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

$\Rightarrow \beta_j = \dfrac{c_j + \lambda}{a_j} < 0 \Rightarrow c_j < -\lambda$   strong neg corr., then $\beta_j < 0$

If $\beta_j > 0$

$a_j\beta_j - c_j + \lambda = 0 \Rightarrow \beta_j = \dfrac{c_j - \lambda}{a_j} \Rightarrow c_j > \lambda$   strong pos. corr. then $\beta_j > 0$

If $\beta_j = 0$    $-\lambda < c_j < \lambda$     if not strong corr., $\beta_j = 0$

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

11

---

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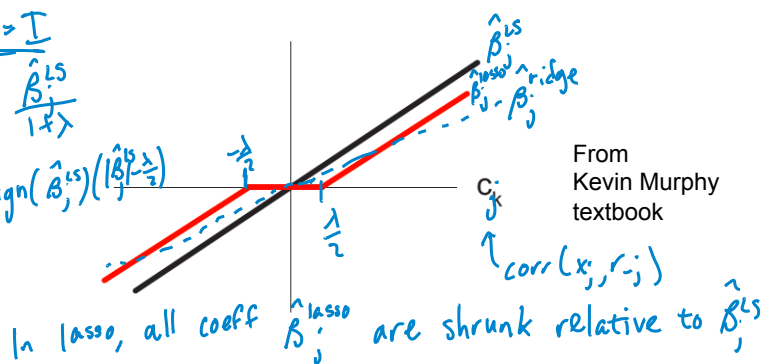
$$\hat{\beta}_j = \begin{cases} (c_j + \lambda)/a_j & c_j < -\lambda \\ 0 & c_j \in [-\lambda, \lambda] \\ (c_j - \lambda)/a_j & c_j > \lambda \end{cases} = \text{sign}\left(\dfrac{c_j}{a_j}\right)\left(\dfrac{|c_j|}{a_j} - \dfrac{\lambda}{a_j}\right)_+$$

If $X^T X = I$

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

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

$\hat{\beta}^{LS}$

$\hat{\beta}^{lasso}$  $\hat{\beta}^{ridge}$

$\hat{\beta}_j - \hat{\beta}_j^{ridge}$

$-\dfrac{\lambda}{2}$

$\dfrac{\lambda}{2}$

$c_k$

From
Kevin Murphy
textbook

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

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

12

6

# Coordinate Descent

- Given a function F$(\beta)$      $\leftarrow F(\beta_1, \ldots, \beta_p)$
  - Want to find minimum    $\beta^* = \min_{\beta} F(\beta)$

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

  *1-d optimization problem... Just solved this for lasso*

- Coordinate descent:

  *While not converged*
  *Pick coord j:*
  $$\beta_j \leftarrow \min_{b} F(\beta_1, \ldots, \beta_{j-1}, b, \beta_{j+1}, \ldots, \beta_p)$$

- How do we pick a coordinate?

  *Round robin, randomly, smartly....*

- When does this converge to optimum?

  *e.g. strongly convex, separability*    *not strongly*
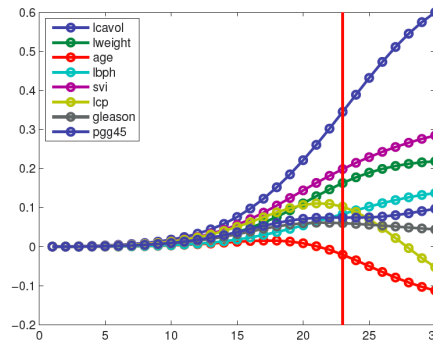
13

---

# 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:   *Cache*
      $$a_j = 2\sum_{i=1}^{N}(x_j^i)^2 \qquad c_j = 2\sum_{i=1}^{N} x_j^i(y^i - \beta'_{-j}x^i_{-j})$$

    - For convergence rates, see Shalev-Shwartz and Tewari 2009

- **Other common technique = LARS**
  - Least angle regression and shrinkage, Efron et al. 2004
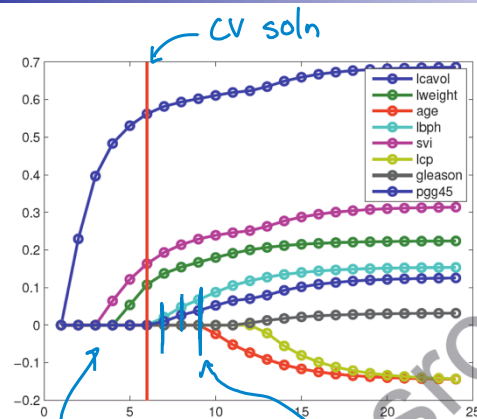
14

7

# Recall: *Ridge Coefficient Path*



From Kevin Murphy textbook

- Typical approach: select λ using cross validation

15

---

# Now: *LASSO Coefficient Path*



From Kevin Murphy textbook

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

*only a few critical values of λ where support changes*

*inc. λ*

*solns are sparse for any given λ*

*CV soln*

$B$

16

# LASSO Example

$\hat{\beta}$ cv solns

| 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 | −0.046 | |
| lbph | 0.210 | 0.162 | 0.002 |
| svi | 0.305 | 0.227 | 0.094 |
| lcp | −0.288 | 0.000 | |
| gleason | −0.021 | 0.040 | |
| pgg45 | 0.267 | 0.133 | |

$\hat{\beta_0}$ $\hat{\beta_1}$ $\hat{\beta_p}$

From Rob Tibshirani slides

not in the model

sparse solns

©Emily Fox 2013                    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

$$\hat{\theta} \rightarrow \theta^* \ ?$$

- 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*

n grows slowly relative to p

©Emily Fox 2013                    19

# 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

  $$1 - 4\exp(-c_1 n \lambda_n^2) \longrightarrow 1:$$

  - The LASSO problem has a unique solution with support contained within the true support $S(\hat{\beta}^{lasso}) \subseteq S(\beta^*)$
  - If $\min_{j \in S(\beta^*)} |\beta_j^*| > c_2 \lambda_n$ for some $c_2 > 0$, then $\boxed{S(\hat{\beta}) = S(\beta^*)}$

20

---

# 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

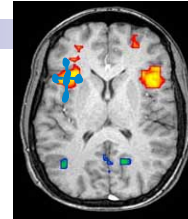  $$\|y - X\beta\|_2^2 + \lambda_1 \sum_j |\beta_j| + \lambda_2 \|\beta\|_2^2$$

  (still some issues, but other solns)

21

10

# 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:

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

*{ in edge set*

*penalizing $\beta_s \neq \beta_t$*

22

# A Bayesian Formulation
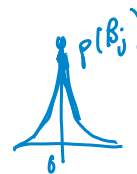
- Consider a model with likelihood

$$y_i \mid \beta \sim N(\beta_0 + x_i^T \beta, \sigma^2)$$

  and prior

$$\beta_j \sim \text{Lap}(\beta_j; \lambda)$$

  where

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

  $p(\beta_j)$

- For large λ    *more peaked around 0*

  $p(\beta_j = 0) = 0$
  *but look at the post. mode*

- LASSO solution is equivalent to the **_mode_** of the posterior
- Note: posterior mode ≠ 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 prior as an alternative*

23

11