

Case Study 3: fMRI Prediction

LASSO Regression

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

Emily Fox

February 19th, 2013

©Emily Fox 2013

1

fMRI Prediction Task

- **Goal:** Predict word stimulus from fMRI image

Can we read your brain?



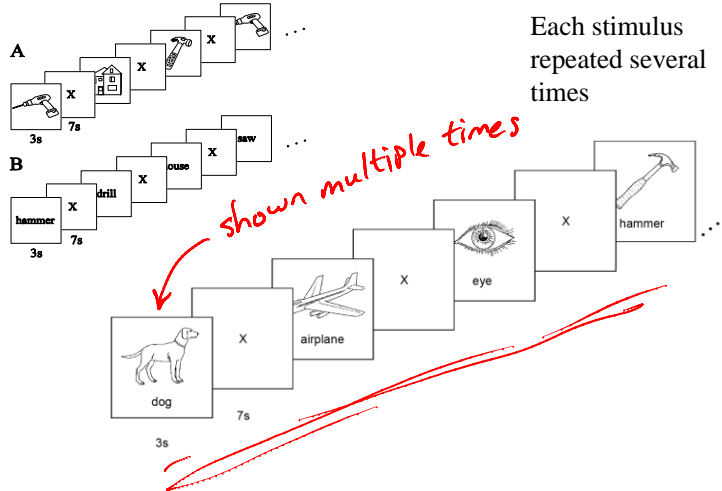
Classifier
(logistic regression,
kNN, ...)

~~HAMMER~~
or
HOUSE

©Emily Fox 2013

2

Typical Stimuli



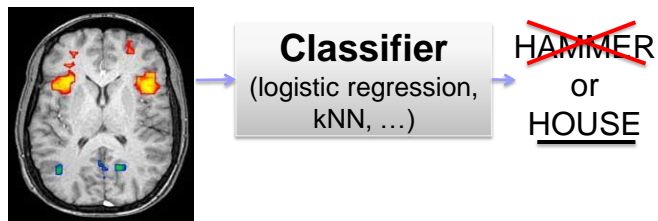
©Emily Fox 2013

3

Zero-Shot Classification

- **Goal:** Classify words not in the training set
- **Challenges:**
 - Cost of fMRI recordings is high
 - Can't get recordings for every word in the vocabulary

Never showed "giraffe" in scanner



©Emily Fox 2013

4

Semantic Features

Google Trillion word corpus

Semantic feature values: "celery"

0.8368, eat
0.3461, taste
0.3153, fill
0.2430, see
0.1145, clean
0.0600, open
0.0586, smell
0.0286, touch
...
...
0.0000, drive
0.0000, wear
0.0000, lift
0.0000, break
0.0000, ride

co-occurrence

Semantic feature values: "airplane"

0.8673, ride
0.2891, see
0.2851, say
0.1689, near
0.1228, open
0.0883, hear
0.0771, run
0.0749, lift
...
...
0.0049, smell
0.0010, wear
0.0000, taste
0.0000, rub
0.0000, manipulate

©Emily Fox 2013

5

Zero-Shot Classification

- From training data, learn two mappings:
 - $A = \{ \text{few} \} \rightarrow \text{"dog"} \}$
 - $B = \{ \text{many} \} \rightarrow \text{"dog"} \}$
- Can use "cheap" co-occurrence data to help learn L
 - Training = $\{ \text{dog} \rightarrow [\cdot] \}$ (use both A + B)
 - N examples ... N small



Predict

Features of word

Classifier
(logistic regression, kNN, ...)

~~HAMMER~~
or
HOUSE

new image $\rightarrow [\cdot] \xrightarrow{\text{using B}} \text{"giraffe"}$

$\rightarrow S \leftarrow \text{learned from training data}$

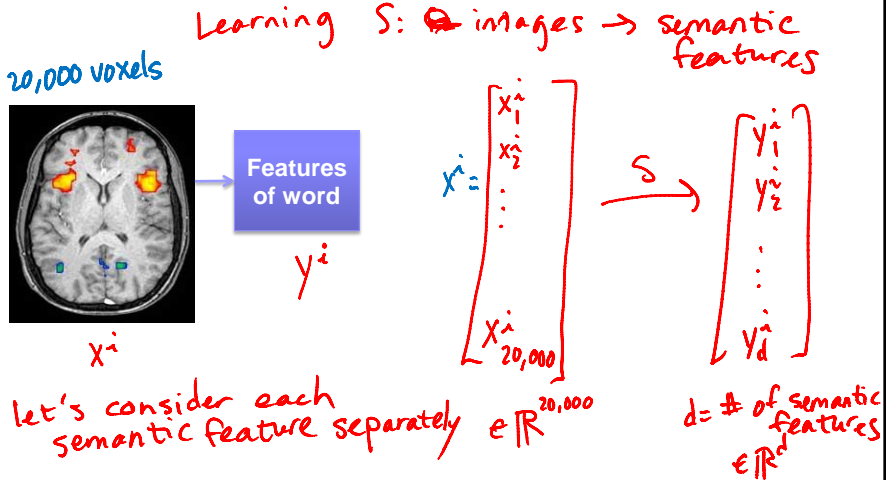
©Emily Fox 2013

6

fMRI Prediction Subtask

Challenge:
 $p \gg N$
 scenario

- Goal: Predict semantic features from fMRI image



©Emily Fox 2013

7

Ridge Regression

- Ameliorating issues with overfitting: penalization of weights = "regularization"
- New objective: \leftarrow same as in LS

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

RSS $\beta^T \beta$

don't penalize intercept term

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

- Reformulate:

$$F(\beta) = \frac{1}{2} \beta^T (X^T X) \beta - \beta^T (X^T y) + \text{const.} + \frac{1}{2} \lambda \beta^T \beta$$

RSS(β)

$$= \frac{1}{2} \beta^T (X^T X + \lambda I) \beta - \beta^T (X^T y) + \text{const.}$$

- Set gradient = 0

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} (X^T y)$$

©Emily Fox 2013

8

Variable Selection

- Ridge regression: Penalizes large weights
- What if we want to perform "feature selection"?
 - 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
 - Stepwise procedures are sensitive to data perturbations and often include features with negligible improvement in fit
- Try new penalty: Penalize non-zero weights
 - Penalty:
 - Leads to sparse solutions
 - Just like ridge regression, solution is indexed by a continuous param λ

Min. this obj. / coeff. are very sensitive to what's inc. in model

discrete 2^P subsets of predictors ... can't do this

greedy, but w/ backtracking. --

$$\|B\|_1 = \sum |B_j| \quad l_1\text{-reg.}$$

LASSO Regression

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

$$\min_B \sum_{i=1}^N (y^i - (b_0 + B^T x^i))^2 + \lambda \|B\|_1$$

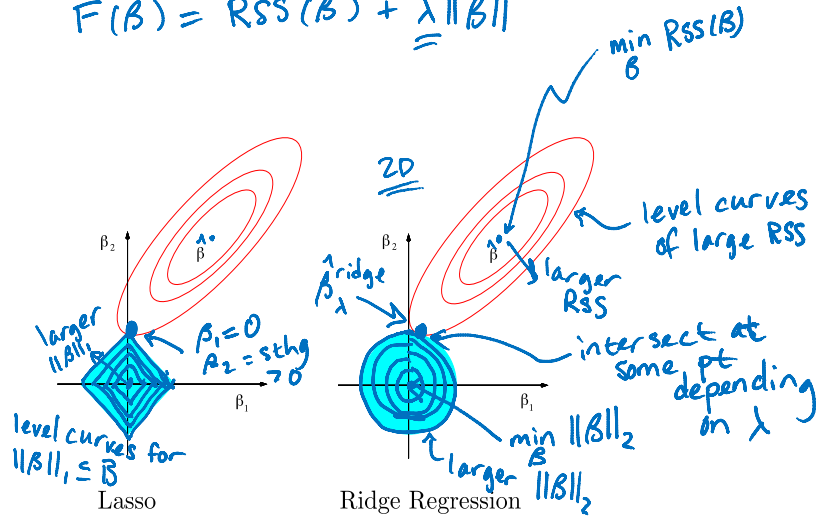
RSS(B)



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

Geometric Intuition for Sparsity

$$F(\beta) = \text{RSS}(\beta) + \lambda \|\beta\|$$



©Emily Fox 2013

11

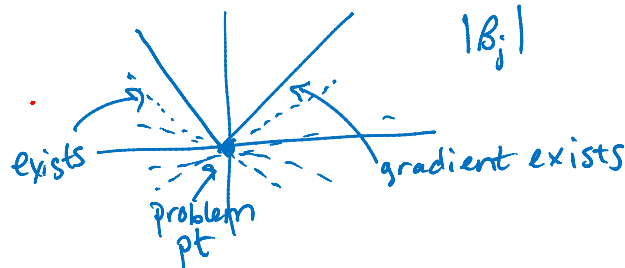
Soft Thresholding

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

look at β_j fixing all others

- L1 penalty $\|\beta\|_1$ is not differentiable whenever $\beta_j = 0$

- Look at subgradient...



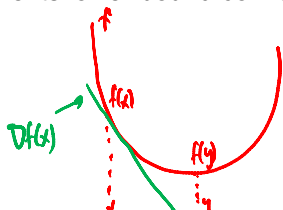
©Emily Fox 2013

12

Subgradients of Convex Functions

From Case Study 1:

- Gradients lower bound convex functions:

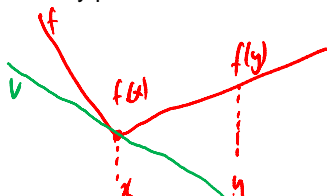


$$f(y) \geq f(x) + Df(x)(y-x)$$

- Gradients are unique at x if function differentiable at x

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



$$v \in \partial f(x) \text{ subgradient}$$

$$\text{if}$$

$$f(y) \geq f(x) + v \cdot (y-x)$$

©Carlos Guestrin 2013

13

Soft Thresholding

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

- Gradient of RSS term:

$$\frac{d}{d\beta_j} \text{RSS}(\beta) = a_j \beta_j - c_j \leftarrow 2 \sum_{i=1}^N x_j^i (y^i - \beta_j^T x_{-j}^i)$$

- Subgradient of full objective:

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

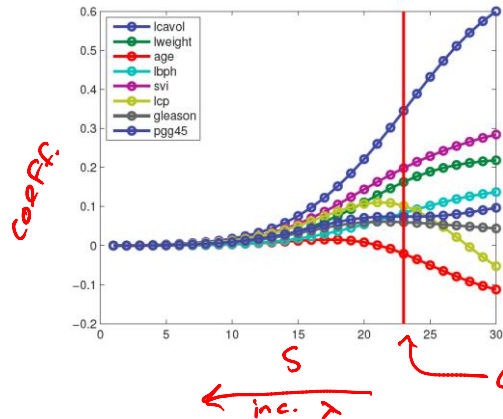
$$= \begin{cases} a_j \beta_j - c_j - \lambda & \beta_j < 0 \\ [-c_j - \lambda, -c_j + \lambda] & \beta_j = 0 \\ a_j \beta_j - c_j + \lambda & \beta_j > 0 \end{cases}$$

$c_j \propto \text{corr}(x_j, \hat{y}_{-j})$
 \hat{y}_{-j} is residual from model w/o j th cov.
 \hat{y}_{-j} is all but the j th coeff.
 λ is msc how relevant x_j is for pred y beyond what the others can

©Emily Fox 2013

14

Recall: Ridge Coefficient Path



From Kevin Murphy textbook

$$\|B\|_2 \leq S$$

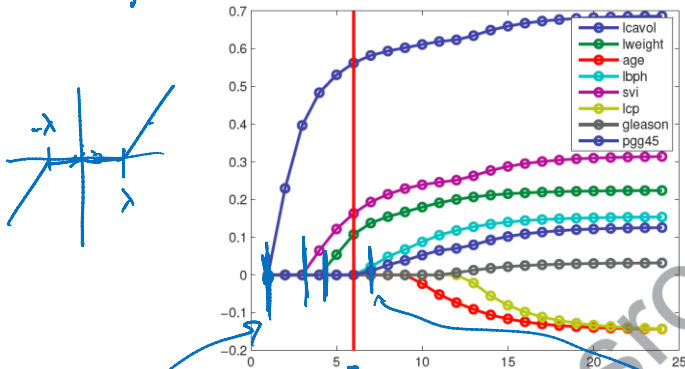
- Typical approach: select λ using cross validation (CV)

©Emily Fox 2013

17

Now: LASSO Coefficient Path

Again, each λ indexes a diff. soln



From Kevin Murphy textbook

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

inc. λ

©Emily Fox 2013

18

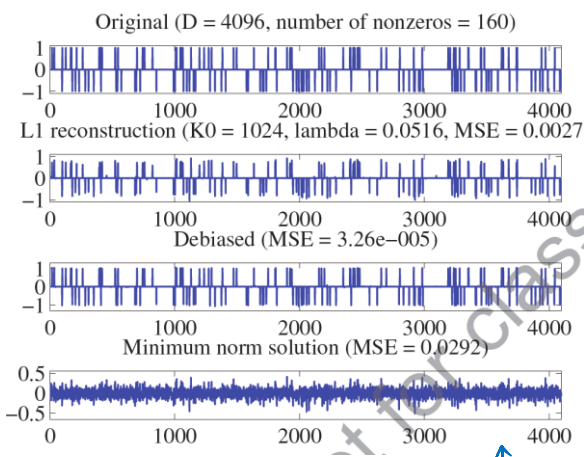
LASSO Example

	Term	<u>Least Squares</u>	Ridge	Lasso
$\hat{\beta}_0$	Intercept	2.465	2.452	2.468
$\hat{\beta}_1$	lcavol	0.680	0.420	0.533
\vdots	lweight	0.263	0.238	0.169
\vdots	age	-0.141	<u>-0.046</u>	
\vdots	lbph	0.210	<u>0.162</u>	0.002
\vdots	svi	0.305	<u>0.227</u>	0.094
\vdots	lcp	-0.288	<u>0.000</u>	
$\hat{\beta}_p$	gleason	-0.021	0.040	
	pgg45	0.267	0.133	

red line (CV) solns

not in model

Debiasing



all coeff. shrunk \rightarrow bias

Some people:

1. Use LASSO to find support
2. Run regression just w/ the selected cov.

\Rightarrow removes bias for this model

From Kevin Murphy textbook

LS est.

Sparsistency

- Typical Statistical Consistency Analysis:

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

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

- 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

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

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

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

Acknowledgements



- Some material in this lecture was based on slides provided by:
 - Tom Mitchell – fMRI
 - Rob Tibshirani – LASSO
 - Ryan Tibshirani – Fused LASSO