

## Case Study 3: fMRI Prediction

# fMRI Prediction Task, Ridge, LASSO Review Fused LASSO, LARS

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

*so far large/streaming N,  
now big-p domain*

Emily Fox  
January 30<sup>th</sup>, 2014

©Emily Fox 2014

1

## fMRI Prediction Task

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

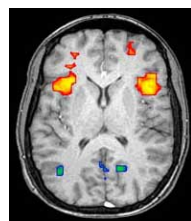
- **Challenges:**

- $p \gg N$  (feature dimension  $\gg$  sample size)
- Cost of fMRI recordings is high
- Only have a few training examples for each word

*# of voxels = # params*

*many more  
params than obs.  
What can we do?*

*20,000 voxels  
(big P)*



**Classifier**

(logistic regression,  
kNN, ...)

~~HAMMER~~

or  
HOUSE

*few training example  
# = N*

©Emily Fox 2014

2

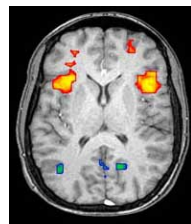
# Zero-Shot Classification

- **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*



**Classifier**  
(logistic regression,  
kNN, ...)

~~HAMMER~~  
or  
HOUSE

©Emily Fox 2014

3

# Semantic Features

*side info =*

*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  
...

*CO-OCCURRENCE*

0.0000, drive  
0.0000, wear  
0.0000, lift  
0.0000, break  
0.0000, ride

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 2014

4

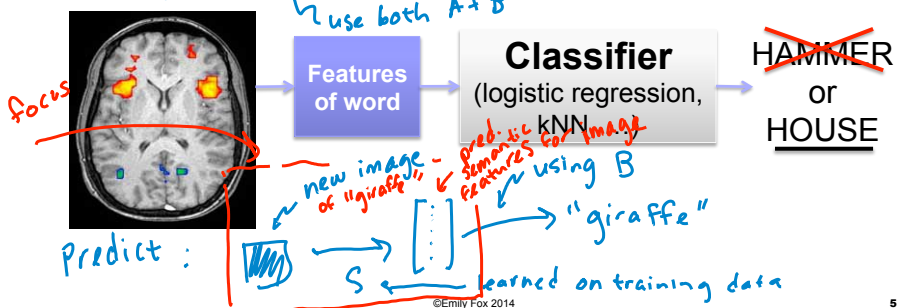
# Zero-Shot Classification

- From training data, learn two mappings

- S: input image  $\rightarrow$  semantic features
- L: semantic features  $\rightarrow$  word

- Can use "cheap" co-occurrence data to help learn L

Training:  $\{ \text{image} \rightarrow [ \cdot ] \rightarrow \text{"dog"} \}$  N examples... N small



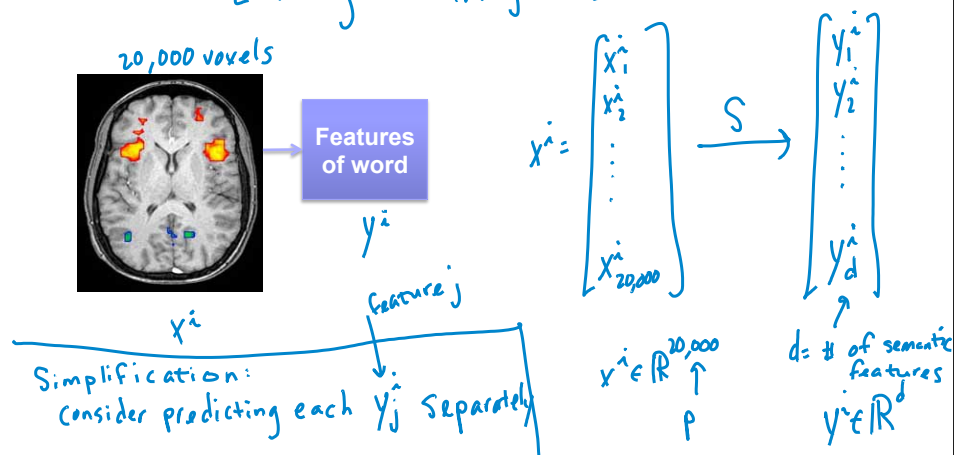
©Emily Fox 2014

5

# fMRI Prediction Subtask

- Goal:** Predict semantic features from fMRI image

Learning  $S$ : images  $\rightarrow$  semantic features



©Emily Fox 2014

6

# Linear Regression

Simplest model

Note: previously, we  $\beta$  more common in stats

Model:  $y^i = \beta_0 + \beta_1 x_1^i + \dots + \beta_p x_p^i + \epsilon^i$   
 one feat.  $\epsilon \in \mathbb{R} \rightarrow = \beta^T x^i + \epsilon^i$

$$\epsilon^i \sim N(0, \sigma^2)$$

$$\Downarrow$$

$$y^i \sim N(\beta^T x^i, \sigma^2)$$

All obs: 
$$\begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix} = \begin{bmatrix} x_1^1 & \dots & x_p^1 \\ \vdots & & \vdots \\ x_1^N & \dots & x_p^N \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \vdots \\ \epsilon^N \end{bmatrix}$$

MLE:  $\hat{\theta} = \arg \max_{\theta} \log p(D | \theta) \leftarrow \sum_{i=1}^N \log p(y^i | x^i, \theta) = \frac{-1}{2\sigma^2} \sum_{i=1}^N (y^i - \beta^T x^i)^2 + \text{const}$

all obs.  $\uparrow$  ind. obs.

RSS( $\beta$ )

$$\hat{\theta}^{ML} = \arg \min_{\beta} \text{RSS}(\beta) = \underbrace{(X^T X)^{-1}}_{(p+1) \times (p+1)} X^T y$$

Here,  $p \gg N$   
 so  $X^T X$  low rank  
 + we want its inverse...

- Minimizing RSS= least squares regression

# Ridge Regression

- Ameliorating issues with overfitting:

penalization of weights = "regularization"

- New objective:

$$\min_{\beta} \sum_{i=1}^N (y^i - (\beta_0 + \beta^T x^i))^2 + \lambda \|\beta\|_2^2$$

don't want to penalize intercept  $\uparrow$

redefine  $X$  w/o 1's  $\leftarrow \beta^T \beta$  where  $\beta = [\beta_1 \dots \beta_p]$

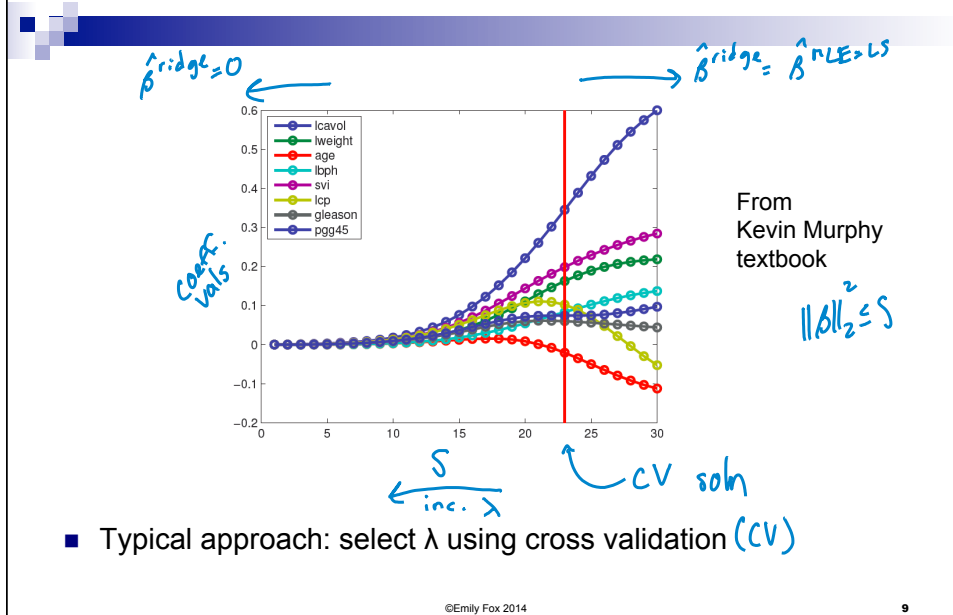
$$\min_{\beta} \text{RSS}(\beta) \quad \text{s.t.} \quad \|\beta\|_2^2 \leq S$$

- Solution:

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

new term  $\uparrow$

# Ridge Coefficient Path



## Case Study 3: fMRI Prediction

### fMRI Prediction Results

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

Emily Fox  
January 30<sup>th</sup>, 2014

# fMRI Prediction Results

- Palatucci et al., "Zero-Shot Learning with Semantic Output Codes", NIPS 2009
- fMRI dataset:
  - 9 participants
  - 60 words (e.g., bear, dog, cat, truck, car, train, ...)
  - 6 scans per word
  - Preprocess by creating 1 "time-average" image per word
- Knowledge bases
  - Corpus5000 – semantic co-occurrence features with 5000 most frequent words in Google Trillion Word Corpus
  - human218 – Mechanical Turk (Amazon.com) 218 semantic features ("is it manmade?", "can you hold it?") Scale of 1 to 5

2 diff. sources of side info  
→ compare performance using each

# fMRI Prediction Results

- **First stage:** Learn mapping from images to semantic features

- Ridge regression

$X \in \mathbb{R}^{N \times p}$  ← fMRI images  
matrix of  $x^i$

$F \in \mathbb{R}^{N \times d}$  ← obs. matrix of  $f^i$   
# semantic features

From (limited) training data:  $\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T F$  ← right soln assuming ind. features

pred. semantic features of new image:  $\hat{f}^{\text{new}} = X^{\text{new}} \hat{\beta}_{\text{ridge}}$

... d ind. subproblems stacked up

- **Second stage:** 1-NN classification using knowledge base

look for word in B w/ f closest to  $\hat{f}^{\text{new}}$

# fMRI Prediction Results

- Leave-two-out-cross-validation

- Learn ridge coefficients using 58 fMRI images
- Predict semantic features of 1<sup>st</sup> heldout image
- Compare whether semantic features of 1<sup>st</sup> or 2<sup>nd</sup> heldout image are closer

Table 1: Percent accuracies for leave-two-out-cross-validation for 9 fMRI participants (labeled P1-P9). The values represent classifier percentage accuracy over 3,540 trials when discriminating between two fMRI images, both of which were omitted from the training set.

	P1	P2	P3	P4	P5	P6	P7	P8	P9	Mean
corpus5000	79.6	67.0	69.5	56.2	77.7	65.5	71.2	72.9	67.9	69.7
human218	90.3	82.9	86.6	71.9	89.5	75.3	78.0	77.7	76.2	80.9

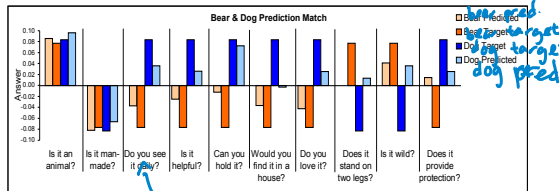


Figure 1: Ten semantic features from the human218 knowledge base for the words bear and dog. The true encoding is shown along with the predicted encoding when fMRI images for bear and dog were left out of the training set.

©Emily Fox 2014

13

*out of 60*

*9 subjects*

*stat. sig.*

*mean target dog pred*

*"yes" for dogs  
"no" for bears  
+ our pred. agree*

*"do you see it every day"*

# fMRI Prediction Results

- Leave-one-out-cross-validation

- Learn ridge coefficients using 59 fMRI images
- Predict semantic features of heldout image
- Compare against very large set of possible other words

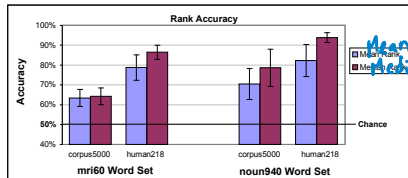


Figure 2: The mean and median rank accuracies across nine participants for two different semantic feature sets. Both the original 60 fMRI words and a set of 940 nouns were considered.

Table 2: The top five predicted words for a novel fMRI image taken for the word in bold (all fMRI images taken from participant P1). The number in the parentheses contains the rank of the correct word selected from 941 concrete nouns in English.

Bear	Foot	Screwdriver	Train	Truck	Celery	House	Pants
(1)	(1)	(1)	(1)	(2)	(5)	(6)	(21)
<b>bear</b>	foot	screwdriver	train	jeep	beet	supermarket	clothing
fox	feet	pin	jet	truck	artichoke	hotel	vest
wolf	ankle	nail	jail	minivan	grape	theater	t-shirt
yak	knee	wrench	factory	bus	cabbage	school	clothes
gorilla	face	dagger	bus	sedan	celery	factory	panties

©Emily Fox 2014

14

*Mean rank Median rank*

*How high did true word fall on the list of ranked words from pred.*

## Case Study 3: fMRI Prediction

# LASSO Review

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

Emily Fox

January 30<sup>th</sup>, 2014

©Emily Fox 2014

15

## 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:  $\|B\|_1 = \sum_j |B_j|$   $L_1$ -reg.
  - Leads to sparse solutions
  - Just like ridge regression, solution is indexed by a continuous param  $\lambda$

*not min. this obj.  
coeff are very  
sensitive to  
what was  
included in  
model*

*discrete*

*$2^p$  subsets of predictors... clearly not feasible for large  $p$*

*← greed, but  $\exists$  backtracking  
approaches*

©Emily Fox 2014

16



# LASSO Regression

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

- New objective:

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



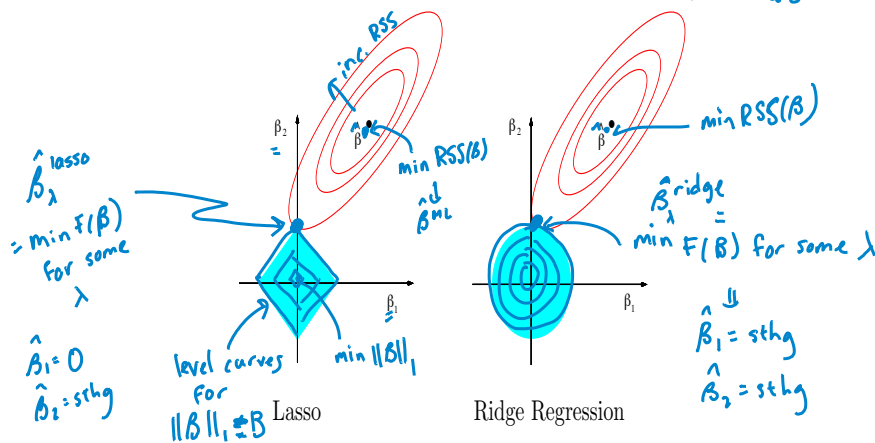
$$\min_{\beta} \text{RSS}(\beta) \quad \text{s.t.} \quad \|\beta\|_1 \leq B$$

©Emily Fox 2014

17

# Geometric Intuition for Sparsity

overall:  $F(\beta) = \text{RSS}(\beta) + \lambda \|\beta\|_1$  ← 1 or 2 norm "lasso" "ridge"



©Emily Fox 2014

18

# Soft Thresholding

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

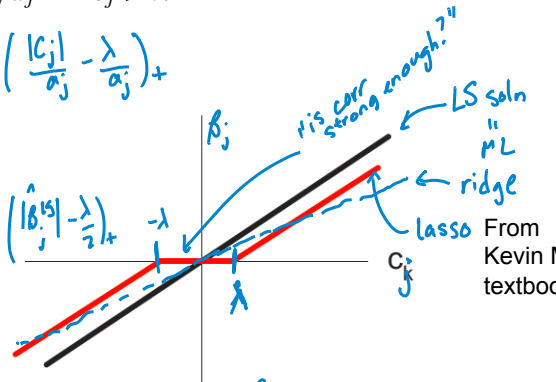
$$= \text{sign}\left(\frac{c_j}{a_j}\right) \left( \frac{|c_j| - \lambda}{a_j} \right)_+$$

If  $X^T X = I$

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

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

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



all examples of feature  $j$   
residual from model w/o using  $j$ th covariate

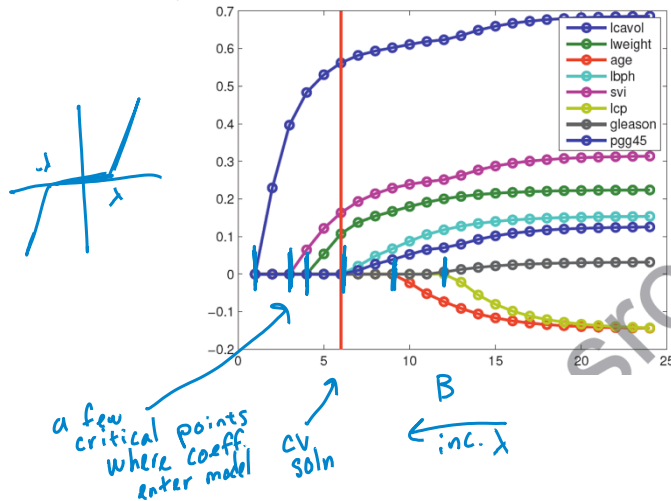
From Kevin Murphy textbook

©Emily Fox 2014

19

# LASSO Coefficient Path

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



From Kevin Murphy textbook

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

©Emily Fox 2014

20

# LASSO Example

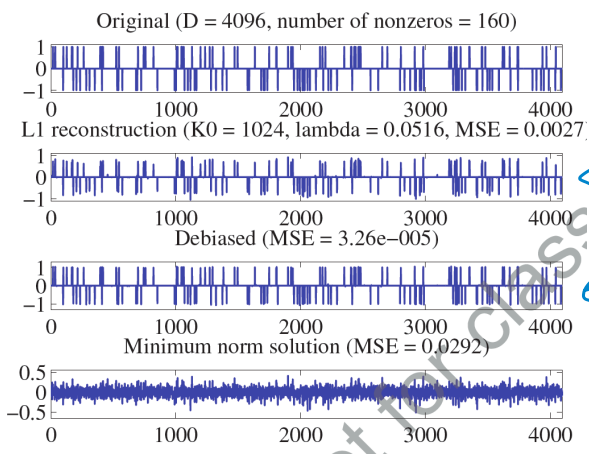
	Term	Least Squares	Ridge	Lasso
$\hat{\beta}_0$	Intercept	2.465	2.452	2.468
$\hat{\beta}_1$	lcavol	0.680	0.420	0.533
$\vdots$	lweight	0.263	0.238	0.169
$\vdots$	age	-0.141	<u>-0.046</u>	
	lbph	0.210	0.162	0.002
	svi	0.305	0.227	0.094
	lcp	-0.288	<u>0.000</u>	
	gleason	-0.021	<u>0.040</u>	
$\hat{\beta}_p$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

CV solns (red lines)

not in the model

shrunk, but non-zero

# Debiasing



all coeff. shrunk → bias

Some people:

1. Use LASSO to find support

2. Run regression just w/ selected covariates

⇒ removes bias for this reduced model

From Kevin Murphy textbook

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

# Sparsistency

- Rescale LASSO objective by  $N$ :
- Theorem (Wainwright 2008, Zhao and Yu 2006, ...):
  - Under some constraints on the design matrix  $X$ , if we solve the LASSO regression using

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

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

# LASSO Algorithms

- Standard convex optimizer
- Least angle regression (LAR)
  - Efron et al. 2004
  - Computes entire path of solutions
  - State-of-the-art until 2008
- Pathwise coordinate descent – new
- More on these “shooting” algorithms next time...

# 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

# Acknowledgements

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