

Case Study 3: fMRI Prediction

fMRI Prediction Task

Now:
 $d \rightarrow p$
 ↑
 dim of
 features/
 Predictors/
 Covariates

Machine Learning for Big Data
 CSE547/STAT548, University of Washington
 Emily Fox
 April 23rd, 2015

"big data challenge":
 large N , streaming N ,
 now big- p domain

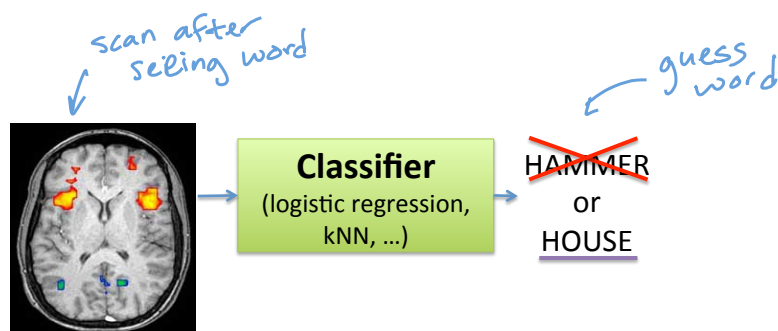
©Emily Fox 2015

1

fMRI Prediction Task

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

can we read your brain?



©Emily Fox 2015

2

fMRI



©Emily Fox 2015

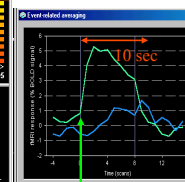
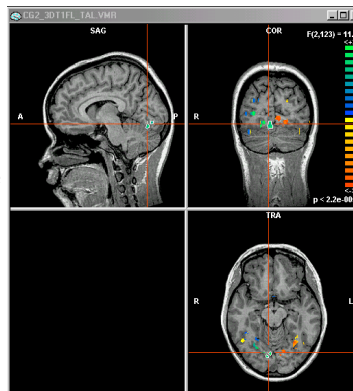
3

fMRI

*high-res**~1 mm resolution
pretty slow
~1 image per sec.*20,000 voxels/image

safe, non-invasive

measures Blood
Oxygen Level
Dependent (BOLD)
response

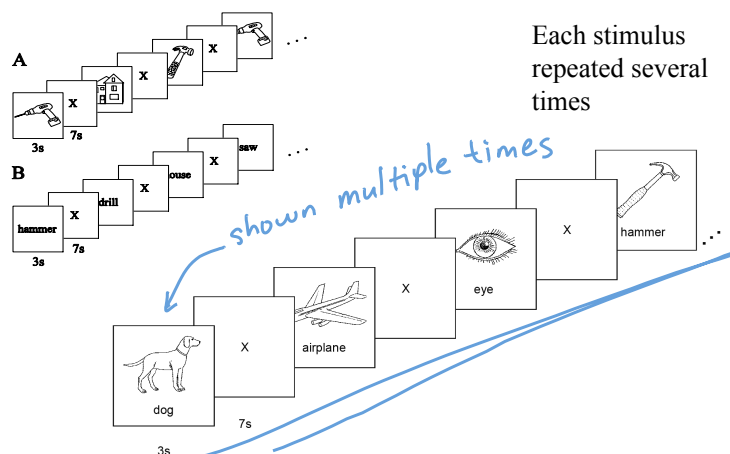


Typical fMRI
response to
impulse of
neural activity

©Emily Fox 2015

4

Typical Stimuli

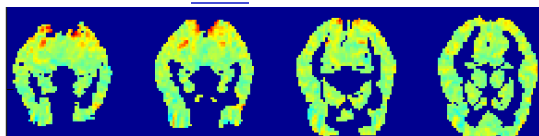


©Emily Fox 2015

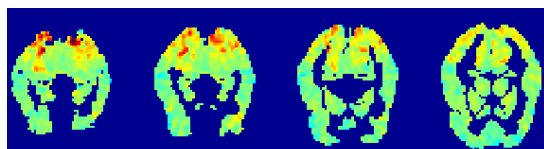
5

fMRI Activation

fMRI activation for "bottle":

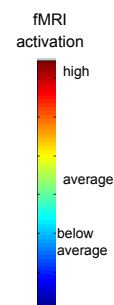
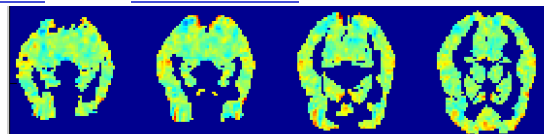


Mean activation averaged over 60 different stimuli:



"bottle" minus mean activation:

is this enough?



©Emily Fox 2015

6

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 = # of params

of obs. of each word

*many more
params than
obs.*

what can we do?



Classifier
(logistic regression,
kNN, ...)

~~HAMMER~~
or
HOUSE

©Emily Fox 2015

7

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

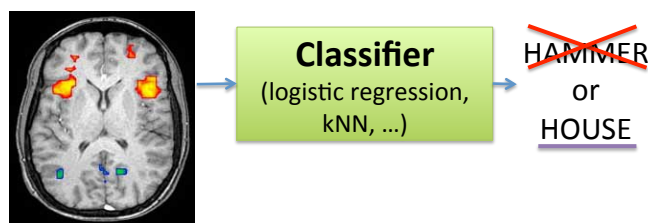
©Emily Fox 2015

8

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
- We don't have many brain images, but we have a lot of info about the words and how they relate (co-occurrence, etc.)
- How do we utilize this "cheap" information?

Many docs that contain "giraffe" also contain "neck" "animal" "zoo" ...



©Emily Fox 2015

9

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 2015

10

Zero-Shot Classification

- From training data, learn two mappings:

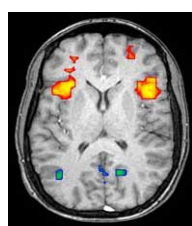
□ S: input image \rightarrow semantic features

□ L: semantic features \rightarrow word

image
 $A = \{ \text{image} \rightarrow \text{"dog"} \}$
 few examples
 Semantic Features
 $B = \{ \text{semantic features} \rightarrow \text{"dog"} \}$
 many examples

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

Trainings $\{ \text{image} \rightarrow \text{semantic features} \rightarrow \text{"dog"} \}$ N examples, N small
 uses both A+B



Features of word

Classifier
 (logistic regression, kNN, ...)

~~HAMMER~~
 or
 HOUSE

Predict:

new image
 $\text{image} \rightarrow \text{semantic features}$

using B (e.g. w/ NN search)
 $\text{semantic features} \rightarrow \text{"giraffe"}$

learned from training data

©Emily Fox 2015

11

fMRI Prediction Subtask

- Goal:** Predict semantic features from fMRI image

Learning S: images \rightarrow semantic features



x^i

Features of word

y^i

$x^i = \begin{bmatrix} x_1^i \\ \vdots \\ x_{20,000}^i \end{bmatrix} \in \mathbb{R}^{20,000}$

$y^i = \begin{bmatrix} y_1^i \\ \vdots \\ y_d^i \end{bmatrix}$

$d = \# \text{ of Semantic Features}$
 $y^i \in \mathbb{R}^d$

Simplification:

Consider predicting each y_j^i separately

©Emily Fox 2015

12

Case Study 3: fMRI Prediction

Ridge, LASSO Review

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

Emily Fox

April 23rd, 2015

©Emily Fox 2015

13

Linear Regression

Model:

$$y^i = \beta_0 + \beta_1 x_1^i + \dots + \beta_p x_p^i + \epsilon_i$$

one feat. $\in \mathbb{R} \rightarrow$

$$= \beta^T x^i + \epsilon_i$$

Note: previously
we "weight"
but β more
common in stat

All obs:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_p^1 \\ \vdots & \vdots & & \vdots \\ x_1^N & x_2^N & \dots & x_p^N \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

β

$$\epsilon_i \sim N(0, \sigma^2)$$

$$\Downarrow$$

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

MLE: $\hat{\theta} = \arg \max_{\theta} \log p(D | \theta)$

$$\sum_{i=1}^N \log p(y^i | x^i, \theta) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (y^i - \beta^T x^i)^2$$

$\text{RSS}(\beta)$
+ const.

$$\hat{\beta}^{ML} = \arg \min_{\beta} \text{RSS}(\beta) = (X^T X)^{-1} X^T y$$

$(p+1) \times N$ $N \times (p+1)$
 $(p+1) \times (p+1)$

Minimizing RSS= least squares regression

©Emily Fox 2015

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

14

Ridge Regression

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

- New objective:

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

Annotations:
 β_0 : don't want to penalize intercept
 β_1, \dots, β_p : [A, ..., Bp]
 x_i : redefine X w/o 1's
 $\beta^T \beta$: $\|\beta\|_2^2$

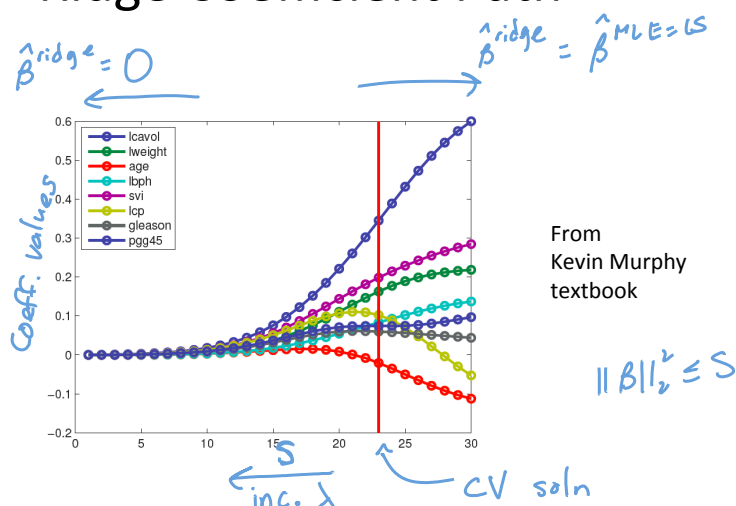
□ Solution: $\min_{\beta} \text{RSS}(\beta) \quad \text{s.t.} \quad \|\beta\|_2^2 \leq S$

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

©Emily Fox 2015

15

Ridge Coefficient Path



From
Kevin Murphy
textbook

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

©Emily Fox 2015

16

Case Study 3: fMRI Prediction

fMRI Prediction Results

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

Emily Fox

April 23rd, 2015

©Emily Fox 2015

17

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 perf. using
each

©Emily Fox 2015

18

fMRI Prediction Results

- First stage: Learn mapping from images to semantic features

- Ridge regression

$$X \in \mathbb{R}^{N \times p}$$
 (matrix of x_i)

$$F \in \mathbb{R}^{N \times d}$$
 (matrix of f_i)

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T F$$
 (right soln assuming d ind. features)

$$\hat{f}^{\text{new}} = X^{\text{new}} \hat{\beta}_{\text{ridge}}$$
 (Subproblems stacked up)

From limited training data: $\hat{\beta}_{\text{ridge}}$
 pred semantic features of new image

obs
 # of semantic features

- Second stage: 1-NN classification using knowledge base

©Emily Fox 2015

19

fMRI Prediction Results

- Leave-two-out-cross-validation

- Learn ridge coefficients using 58 fMRI images
- Predict semantic features of 1st heldout image
- Compare whether semantic features of 1st or 2nd 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.

Google → Mech. Turk →

	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

← stat. sig.

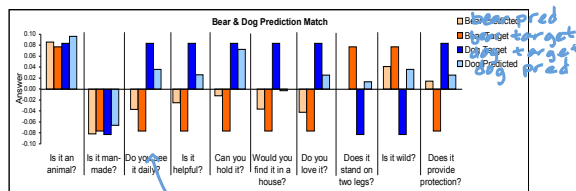


Figure 1: Ten semantic features from the human218 knowledge base for the words bear and dog. The true encoding is shown along with the predicted encoding when fMRI images for bear and dog were left out of the training set.

©Emily Fox 2015

20

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

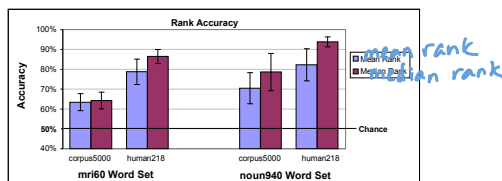


Figure 2: The mean and median rank accuracies across nine participants for two different semantic feature sets. Both the original 60 fMRI words and a set of 940 nouns were considered.

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)
bear	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

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

©Emily Fox 2015

21

Case Study 3: fMRI Prediction

LASSO Review

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

Emily Fox

April 23rd, 2015

©Emily Fox 2015

22

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

not min this obj
coeff. are
very sensitive
to what
was inc. in
the model

discrete

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

- Stepwise procedures are sensitive to data perturbations and often include features with negligible improvement in fit ← greedy, but \exists backtracking approaches

- Try new penalty: Penalize non-zero weights

- Penalty:

$$\|B\|_1 = \sum_j |B_j| \quad L_1\text{-reg.}$$

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

©Emily Fox 2015

23

LASSO Regression

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

- New objective:

$$\min_B \sum_{i=1}^n (y^i - (B_0 + B^T x^i))^2 + \lambda \|B\|_1$$

$\underbrace{\hspace{10em}}_{RSS(B)}$

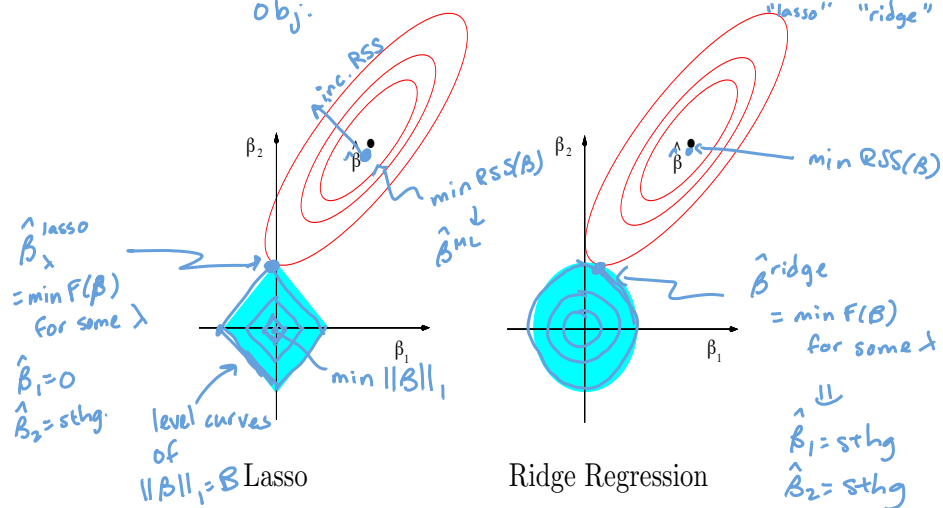
$$\min_B RSS(B) \quad \text{s.t.} \quad \|B\|_1 \leq B$$

©Emily Fox 2015

24

Geometric Intuition for Sparsity

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



©Emily Fox 2015

25

Soft Thresholding

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

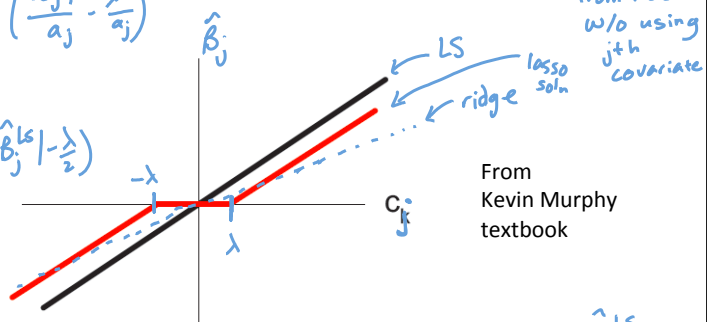
strength of penalty
 $c_j \propto \text{corr}(x_j, r_{-j})$
 all examples of feature j covariate
 residual from model w/o using j th covariate

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

$$\text{If } X^T X = I$$

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

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



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

©Emily Fox 2015

26

Acknowledgements

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