**Case Study 3: fMRI Prediction**

LASSO ~~Regression~~

LARS, Fused LASSO

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

Emily Fox

February 21th, 2013

1

---

# LASSO Regression

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

- New objective:

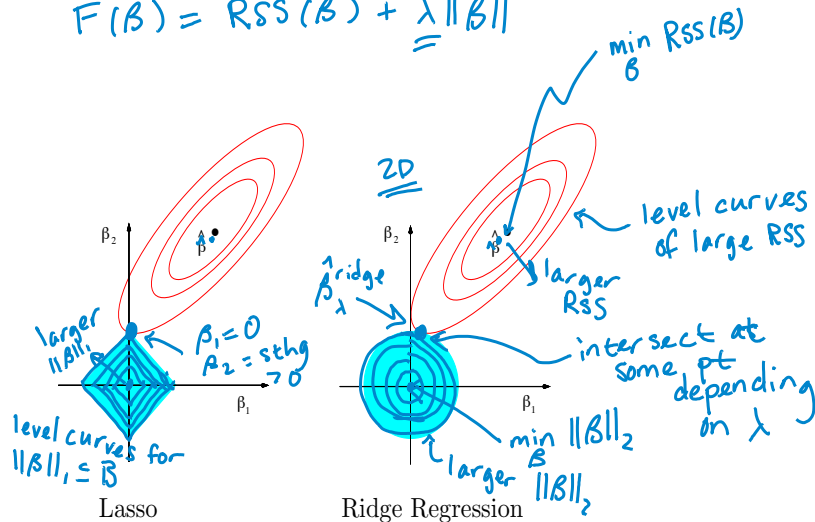$$\min_{\beta} \underbrace{\sum_{i=1}^{N} \left(y^i - (\beta_0 + \beta^T x^i)\right)^2}_{RSS(\beta)} + \lambda \|\beta\|_1$$

$$\Updownarrow$$

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

2

1

# Geometric Intuition for Sparsity
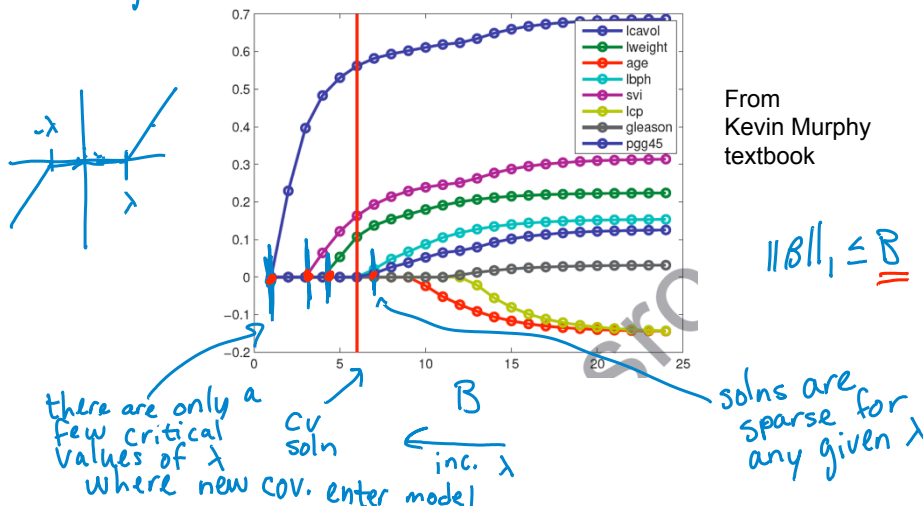
$$F(\beta) = RSS(\beta) + \lambda \|\beta\|$$

$\min\limits_{\beta} RSS(\beta)$

2D

level curves of large RSS

larger RSS

intersect at some pt depending on $\lambda$

$\beta_2$

$\hat{\beta}$

larger $\|\beta\|_1$

$\beta_1 = 0$
$\beta_2 = $ sthg $\neq 0$

$\beta_1$

$\hat{\lambda}^{ridge}$

$\min\limits_{\beta} \|\beta\|_2$

larger $\|\beta\|_2$

level curves for $\|\beta\|_1 \leq \tilde{B}$

Lasso                    Ridge Regression

©Emily Fox 2013                                          3

---

# Now: *LASSO Coefficient Path*

Again, each $\lambda$ indexes a diff. soln



From Kevin Murphy textbook

$\|\beta\|_1 \leq \underline{B}$

there are only a few critical values of $\lambda$ where new cov. enter model

CV soln

$B$

inc. $\lambda$

solns are sparse for any given $\lambda$

©Emily Fox 2013                                          4

2

# 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…

5

# LARS – Efron et al. 2004

- LAR is an efficient stepwise variable selection algorithm
  - "useful and less greedy version of traditional forward selection methods"

  *Efron*

- Can be modified to compute regularization path of LASSO
  - → LARS (Least angle regression and *shrinkage*)

- Increasing upper bound $B$, coefficients gradually "turn on"
  - Few critical values of $B$ where support changes
  - Non-zero coefficients increase or decrease linearly between critical points
  - Can solve for critical values analytically

  *key to providing full reg path*

- Complexity:

$$O(\min(Np^2, pN^2))$$

*# of obs.*    *# of covariates*   = *cost of a single LS soln*

6

3

# LASSO Coefficient Path



From
Kevin Murphy
textbook

---

# LARS – Algorithm

- Assumptions:    *standardize*
    - Response has 0 mean

$$\sum_i y^i = 0$$

    - Covariates are normalized

$$\sum_i x^i_j = 0 \qquad \sum_i (x^i_j)^2 = 1 \qquad j = 1, \ldots, P$$
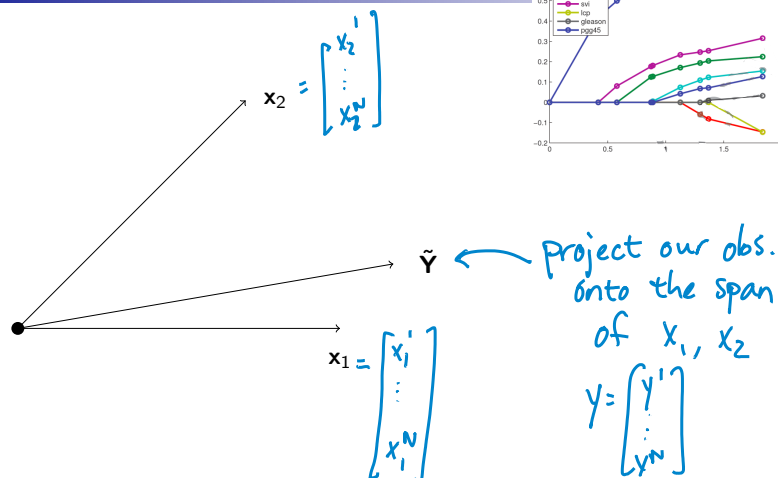
# LARS – Algorithm Overview

- Start with all coefficient estimates $\hat{\beta}_1 = \hat{\beta}_2 = \cdots = \hat{\beta}_p = 0$

- Let $\mathcal{A}$ be the "active set" of covariates most correlated with the "current" residual $\leftarrow$ based on covariates already in model

- Initially, $\mathcal{A} = \{x_{j_1}\}$ for some covariate $x_{j_1}$

- Take the largest possible step in the direction of $x_{j_1}$ until another covariate $x_{j_2}$ enters $\mathcal{A}$

- Continue in the direction equiangular between $x_{j_1}$ and $x_{j_2}$ until a third covariate $x_{j_3}$ enters $\mathcal{A}$

- Continue in the direction equiangular between $x_{j_1}, x_{j_2}, x_{j_3}$ until a fourth covariate $x_{j_4}$ enters $\mathcal{A}$

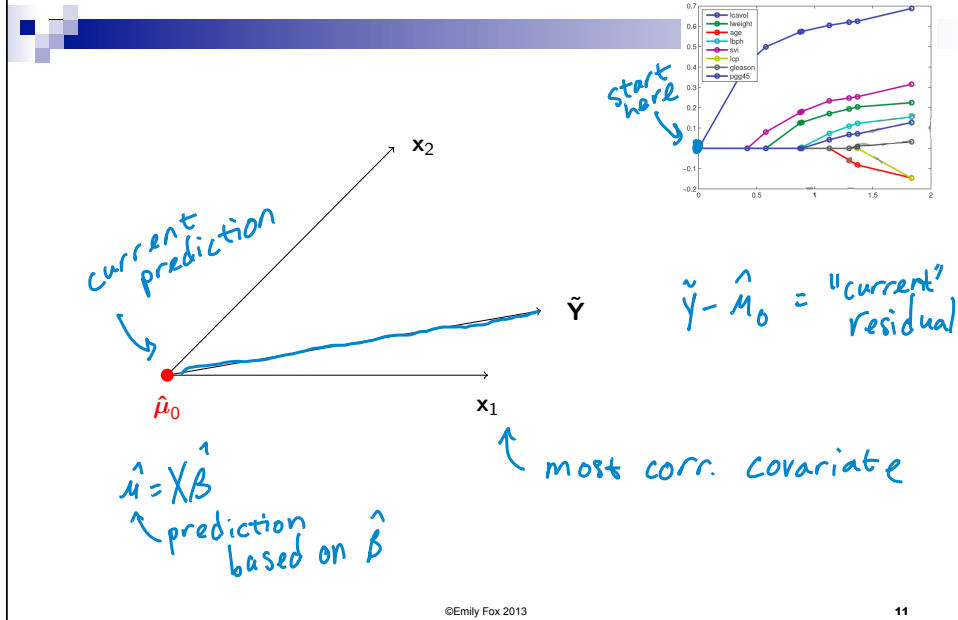- This procedure continues until all covariates are added at which point

9

# LARS – Illustration for *p*=2 covariates



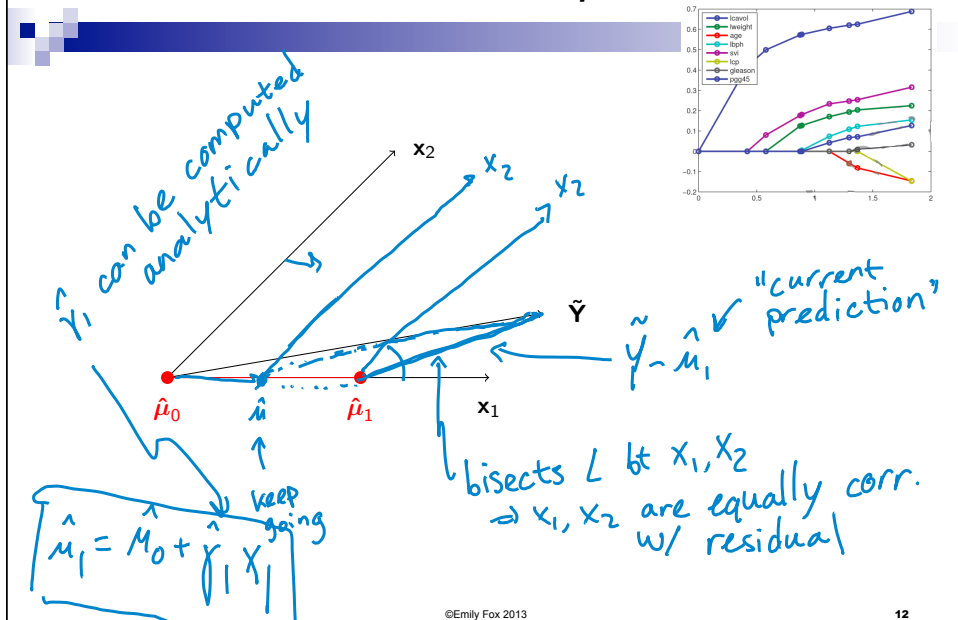$\mathbf{x}_2 = \begin{bmatrix} x_2^1 \\ \vdots \\ x_2^N \end{bmatrix}$

$\tilde{\mathbf{Y}}$ $\leftarrow$ project our obs. onto the span of $x_1, x_2$

$\mathbf{x}_1 = \begin{bmatrix} x_1^1 \\ \vdots \\ x_1^N \end{bmatrix}$

$y = \begin{bmatrix} y^1 \\ \vdots \\ y^N \end{bmatrix}$

10

5

# LARS – Illustration for *p*=2 covariates



current prediction

$\mathbf{x}_2$

$\tilde{\mathbf{Y}}$

$\tilde{y} - \hat{\mu}_0 = $ "current" residual

$\hat{\mu}_0$

$\mathbf{x}_1$

$\hat{\mu} = X\hat{\beta}$

prediction based on $\hat{\beta}$

most corr. covariate

start here

©Emily Fox 2013

11

# LARS – Illustration for *p*=2 covariates



$\hat{\gamma}_1$ can be computed analytically

$\mathbf{x}_2$

$x_2$

$x_2$

$\tilde{\mathbf{Y}}$

$\tilde{y} - \hat{\mu}_1$   "current prediction"

$\hat{\mu}_0$   $\hat{\mu}$   $\hat{\mu}_1$   $\mathbf{x}_1$

keep going

$\hat{\mu}_1 = \hat{\mu}_0 + \hat{\gamma}_1 x_1$

bisects $\angle$ bt $x_1, x_2$
⇒ $x_1, x_2$ are equally corr. w/ residual

©Emily Fox 2013

12

6

# LARS – Illustration for *p*=2 covariates



$x_2$ enters picture

$x_2$

$x_2$

$\tilde{Y}$

$u_2$ $\hat{\mu}_2$

$x_2$ enters

$\hat{\mu}_0$

$\hat{\mu}_1$

$x_1$

walk in dir "equiangular" bt $x_1, x_2$

$$\hat{\mu}_2 = \hat{\mu}_1 + \hat{\gamma}_2 \, u_2$$

13

# LARS – Illustration for *p*=2 covariates



$x_2$

$x_2$

$\tilde{Y} = \hat{\mu}_2$

$\hat{\mu}_0$

$\hat{\mu}_1$

$x_1$

with 2 cov, when both are in the model, $\hat{\mu}_2 = \tilde{Y}$

14

7

# LARS-LASSO Relationship

- Let $\mu(\gamma) = X\beta(\gamma)$ with $\quad \beta_j(\gamma) = \hat{\beta}_j + \gamma \hat{d}_j \quad \leftarrow$ comes from LS soln based on active set

- We showed that for active covariate $j$: $\quad \text{sign}(\hat{\beta}_j) = \text{sign}(\underbrace{x'_j(y - \hat{\mu})}_{c_j})$

$x_2$

$\tilde{Y}$

$\hat{\mu}_0 \qquad \mu(\gamma) \qquad \hat{\mu}_1 = \mu(\hat{\beta}_1) \quad x_1$

incr. $\gamma$

$c_j$

corr. bt $X_j$ and our residual w/o $X_j$ in model

©Emily Fox 2013                                      15

---

# Soft Threshholding

FROM LAST TIME

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

If $X^T X = I$

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

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

$\leftarrow$ LS

ridge

$-\lambda$

$\lambda$

$c_k$

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

$\hat{\beta}_j^{OLS}$

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

From Kevin Murphy textbook

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

©Emily Fox 2013                                      16

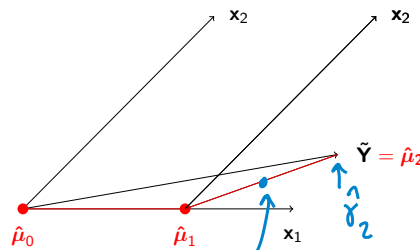# LARS-LASSO Relationship

- Let $\mu(\gamma) = X\beta(\gamma)$ with $\underbrace{\beta_j(\gamma) = \hat{\beta}_j + \gamma \hat{d}_j}$

- We showed that for active covariate *j:* $\quad \text{sign}(\hat{\beta}_j) = \text{sign}(x_j'(y - \hat{\mu}))$

- $\beta_j(\gamma)$ changes sign at $\quad \beta_j(\gamma) = 0 \implies \gamma = -\dfrac{\hat{\beta}_j}{\hat{d}_j}$ $\quad$ Violated

- 1st sign change occurs at $\tilde{\gamma} = \min\limits_{\gamma_j > 0}\{\gamma_j\}$ for covariate $\tilde{j}$

  1st cov. to change signs

---

# LARS-LASSO Relationship

- If $\tilde{\gamma}$ occurs before $\hat{\gamma}$, then next LARS step is not a LASSO solution



$x_2$     $x_2$

$\tilde{Y} = \hat{\mu}_2$

$\hat{\gamma}_2$

$\hat{\mu}_0 \qquad \hat{\mu}_1 \qquad x_1$

$\tilde{\gamma}$ stop here

- **LASSO modification:**

If $\tilde{\gamma} < \hat{\gamma}$, then stop LARS at $\gamma = \tilde{\gamma}$ and remove $\tilde{j}$ from our calc. of equiangular dir.

( remove $\beta_{\tilde{}}$ from model & restart )

# LASSO Coefficient Path



From
Kevin Murphy
textbook

# Comments

- LARS increases $\mathcal{A}$, but LASSO allows it to decrease

- Only involves a single index at a time

- If $p > N$, LASSO returns at most $N$ variables

- If group of variables are highly correlated, LASSO tends to choose one to include rather arbitrarily
  - Straightforward to observe from LARS algorithm….Sensitive to noise.

*beware of interpreting the variables included*

# 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

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

(there some issues... details KM book)

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:

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

penalize these, taking diff vals

has edge in graph

22

11

# Generalized LASSO

- Assume a structured linear regression model:

$$\|y - X\beta\|_2^2 + \lambda \|D\beta\|_1$$

$$D \in \mathbb{R}^{m \times P}$$

- If *D* is invertible, then get a new LASSO problem if we substitute

$$\alpha = D^{-1}\beta$$

- Otherwise, not equivalent

- For solution path, see
Ryan Tibshirani and Jonathan Taylor, "The Solution Path of the Generalized Lasso." Annals of Statistics, 2011.
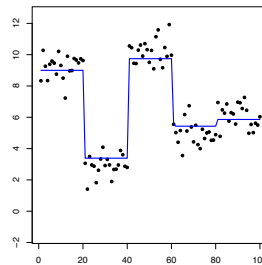
23

---

# Generalized LASSO

signal approximation
scenario
$$X = I$$

$$\hat{\beta}_\lambda = \operatorname*{argmin}_{\beta \in \mathbb{R}^n} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 1 & 0 & 0 & \dots \\ 0 & -1 & 1 & 0 & \dots \\ 0 & 0 & -1 & 1 & \dots \\ \vdots \end{bmatrix}$. This is the 1d fused lasso.

$$\lambda \sum_j |\beta_j - \beta_{j-1}|$$

encourage
piecewise const.
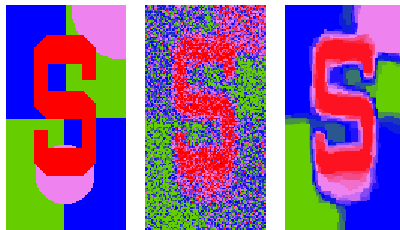~~linear~~

24

12

## Generalized LASSO

$$\hat{\beta}_\lambda = \operatorname*{argmin}_{\beta \in \mathbb{R}^n} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Suppose $D$ gives "adjacent" differences in $\beta$:

$$D_i = (0, 0, \ldots -1, \ldots, 1, \ldots 0),$$

$\lambda \sum_{(s,t) \in E} |\beta_t - \beta_s|$

where adjacency is defined according to a graph $\mathcal{G}$. For a 2d grid, this is the 2d fused lasso.



*encourages constant regions*
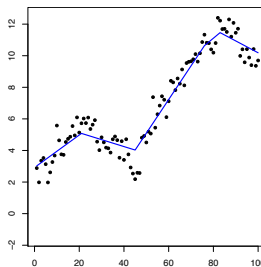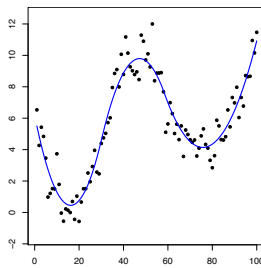
25

## Generalized LASSO

$$\hat{\beta}_\lambda = \operatorname*{argmin}_{\beta \in \mathbb{R}^n} \ \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 2 & -1 & 0 & \ldots \\ 0 & -1 & 2 & -1 & \ldots \\ 0 & 0 & -1 & 2 & \ldots \\ \vdots & & & & \end{bmatrix}$. This is linear trend filtering.



$\beta_3$

$\beta_1 \quad \beta_2$

$\beta_3 - \beta_2 = \beta_2 - \beta_1$

$\Rightarrow \ \{ 2\beta_2 - \beta_1 - \beta_3 = 0$

26

13

# Generalized LASSO

$$\hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \; \frac{1}{2}\|y - \beta\|_2^2 + \lambda\|D\beta\|_1$$

Let $D = \begin{bmatrix} -1 & 3 & -3 & 1 & \dots \\ 0 & -1 & 3 & -3 & \dots \\ 0 & 0 & -1 & 3 & \dots \\ \vdots & & & & \end{bmatrix}$. Get quadratic trend filtering.
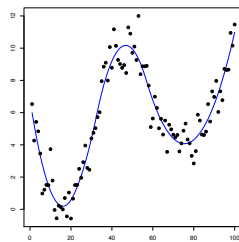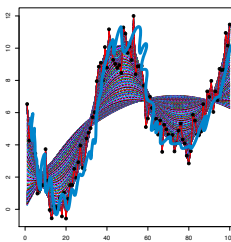
27

---

# Generalized LASSO

- Tracing out the fits as a function of the regularization parameter



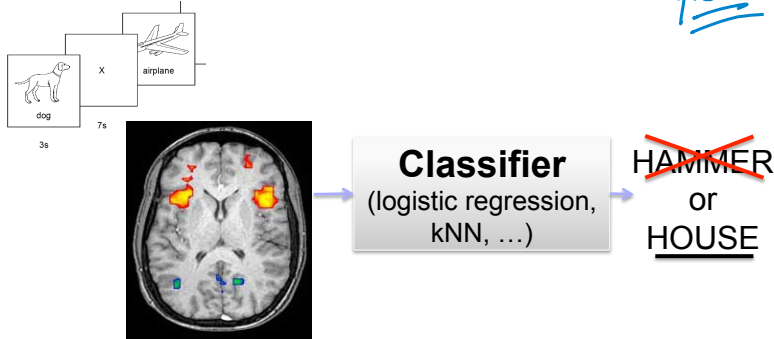$\hat{\beta}_\lambda$ for $\lambda = 25$      $\hat{\beta}_\lambda$ for $\lambda \in [0, \infty]$

28

# fMRI Prediction Task

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

*Can we read your brain?*

*P >> N*
*yes!*



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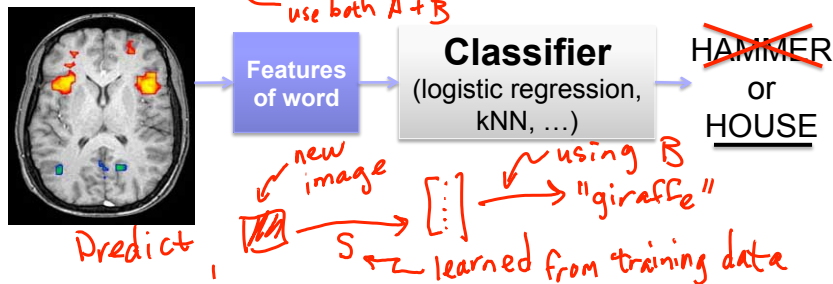(logistic regression,
kNN, …)

~~HAMMER~~
or
HOUSE

---

# Zero-Shot Classification

- From training data, learn two mappings:
  - □ S: input image → semantic features
  - □ L: semantic features → word

*Key →*

$A = \{$ 🖼 $\rightarrow$ "dog"$\}$
*few*

$B = \{[:] \rightarrow$ "dog"$\}$
*many*

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

*from B*

*Training* $= \{$ 🖼 $\rightarrow [:] \rightarrow$ "dog"$\}$   *N examples … N small*

*use both A + B*

**Features of word**

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

~~HAMMER~~
or
HOUSE

*new image*

*Predict* 🖼 $\rightarrow [:] \rightarrow$ "giraffe"

*using B*

*S* ← *learned from training data*

# Semantic Features

*Google Trillion word corpus* (handwritten)

| Semantic feature values: "**celery"** | Semantic feature values: "**airplane"** |
|---|---|
| 0.8368, eat | 0.8673, ride |
| 0.3461, taste | 0.2891, see |
| 0.3153, fill | 0.2851, say |
| 0.2430, see | 0.1689, near |
| 0.1145, clean | 0.1228, open |
| 0.0600, open | 0.0883, hear |
| 0.0586, smell | 0.0771, run |
| 0.0286, touch | 0.0749, lift |
| … | … |
| … | … |
| 0.0000, drive | 0.0049, smell |
| 0.0000, wear | 0.0010, wear |
| 0.0000, lift | 0.0000, taste |
| 0.0000, break | 0.0000, rub |
| 0.0000, ride | 0.0000, manipulate |

*Co-occurrence* (handwritten)

©Emily Fox 2013                                                           31

---

# 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

©Emily Fox 2013                                                           32

16

# fMRI Prediction Results

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

- Ridge regression

$$X \in \mathbb{R}^{N \times P} \rightarrow F \in \mathbb{R}^{N \times d} \leftarrow \text{\# sem. features}$$

*From training data*
$$\hat{\beta}^{ridge} = (X^T X + \lambda I)^{-1} X^T F \leftarrow \text{stack up solns for each sem. feature}$$

$$\hat{f}^{new} = X^{new} \hat{\beta}^{ridge}$$

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

*look for word w/ f closest to $\hat{f}^{new}$*
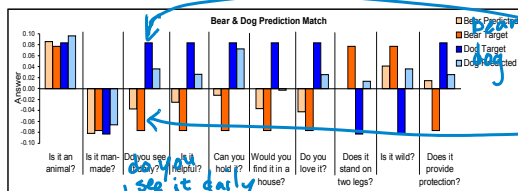
©Emily Fox 2013

33

---

# 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.

| | 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.*

*yes, see dogs daily*

*"bear" ...no*



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 2013

34

17

# fMRI Prediction Results

- Leave-one-out-cross-validation
  - Learn ridge coefficients using 59 fMRI images
  - Predict semantic features of heldout image
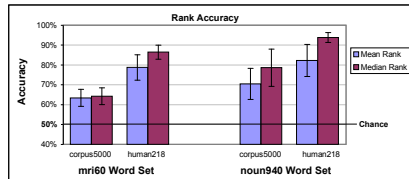  - Compare whether very large set of possible other words



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 |

35

# Acknowledgements

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

36