

Case Study 2: Document Retrieval

Spectral Clustering

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

Emily Fox

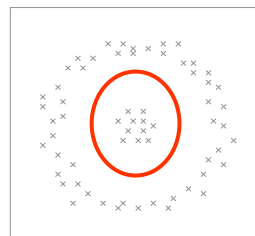
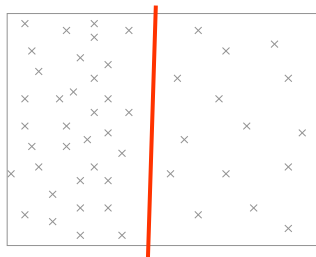
February 14th, 2013

©Emily Fox 2013

1

New Approach: Spectral Clustering

- **Goal:** Cluster observations
- **Method:**
 - Use similarity metric between observations
 - Form a similarity graph
 - Use standard linear algebra and optimization techniques to cut graph into connected components (clusters)



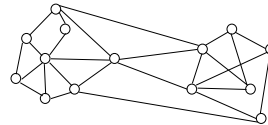
©Emily Fox 2013

2

Setup

- Data: x^1, \dots, x^N
- Similarity metric: s_{ij}

- Similarity graph
 - Nodes v^i
 - Edge weights $w_{ij} = f(s_{ij})$



$$G = \{V, E\}$$

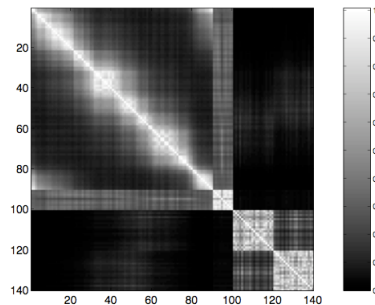
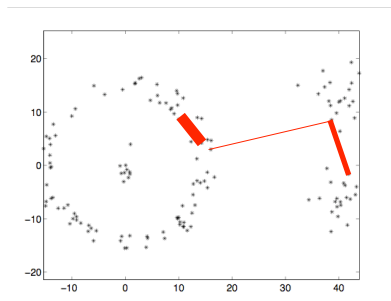
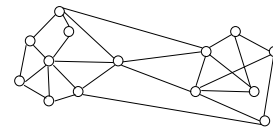
- Problem: Want to partition graph such that edges between groups have low weights

©Emily Fox 2013

3

Graph Terminology I

- Weighted adjacency matrix

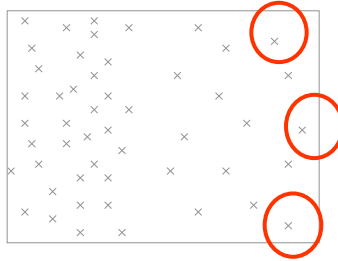


©Emily Fox 2013

4

Issues with MinCut

- MinCut favors isolated clusters

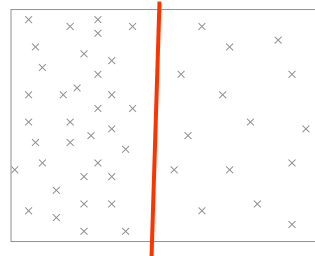


©Emily Fox 2013

5

Cuts Accounting for Size

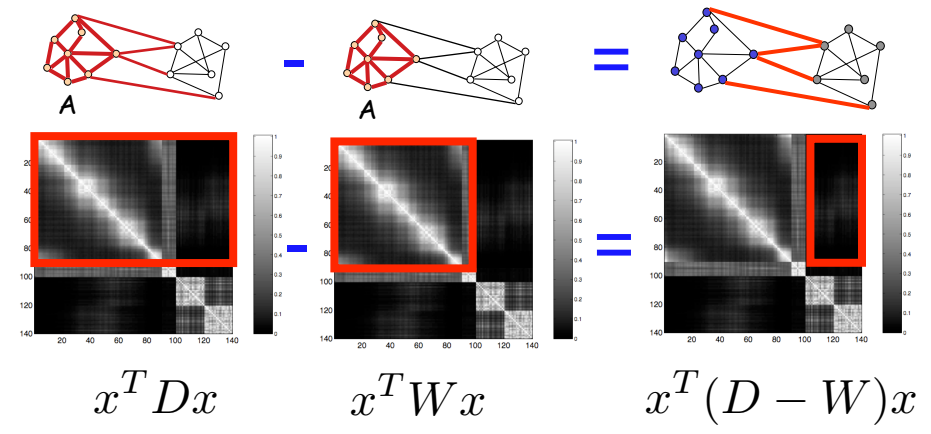
- Ratio cuts (RatioCut)
- Normalized cuts (Ncut)
- Lead to “balanced” clusters



©Emily Fox 2013

6

Restating Cut Metric



©Emily Fox 2013

7

Ratio Cuts for General k

- Define cluster indicator variables:

$$F_{ij} = \begin{cases} 1/\sqrt{|A_j|} & v^i \in A_j \\ 0 & \text{otherwise} \end{cases} \quad \begin{matrix} F'_A F_A = I \\ F_A \in \mathbb{R}^{N \times k} \end{matrix}$$

- RatioCut

$$\text{RatioCut}(A_1, \dots, A_k) = \sum_{i=1}^k f'_{A_i} L f_{A_i} = \text{Tr}(F'_A L F_A)$$

- Reformulating RatioCut problem

$$\min_{A_1, \dots, A_k} \text{Tr}(F'_A L F_A) \quad \text{s.t.} \quad F'_A F_A = I$$

- Relaxation

$$\min_{F \in \mathbb{R}^{N \times k}} \text{Tr}(F' L F) \quad \text{s.t.} \quad F' F = I$$

©Emily Fox 2013

8

Normalized Cuts for General k

- Define cluster indicator variables:

$$F_{ij} = \begin{cases} 1/\sqrt{\text{vol}(A_j)} & v_i \in A_j \\ 0 & \text{ow} \end{cases} \quad \begin{matrix} F'_A F_A = I \\ F'_A D F_A = I \end{matrix}$$

- Reformulating RatioCut problem

$$\min_{A_1, \dots, A_k} \text{Tr}(F'_A L F_A) \quad \text{s.t.} \quad F'_A D F_A = I$$

- Relaxation

$$\min_{H \in R^{N \times k}} \text{Tr}(H' D^{-1/2} L D^{-1/2} H) \quad \text{s.t.} \quad H' H = I$$

- Solution:

- H is matrix of first k eigenvectors of L_{sym} , which is equivalent to the approximate F being the first k eigenvectors of L_{rw}

©Emily Fox 2013

9

Random Walks on Graphs

- Stochastic process with random jumps from v_i to v_j wp:

- Transition matrix:

- Connection to graph Laplacian:

- Intuitively, want to partition graph s.t. random walk stays in cluster for a while and rarely jumps between clusters

©Emily Fox 2013

10

Random Walks on Graphs

- Assume that stationary distribution exists and is unique. Then,
- Proposition: $\text{Ncut}(A, \bar{A}) = P(A | \bar{A}) + P(\bar{A} | A)$
- Proof:
- Minimizing normalized cuts is equivalent to minimizing the probability of transitioning between clusters

©Emily Fox 2013

11

Case Study 3: fMRI Prediction

fMRI Prediction Task, LASSO Regression

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

Emily Fox

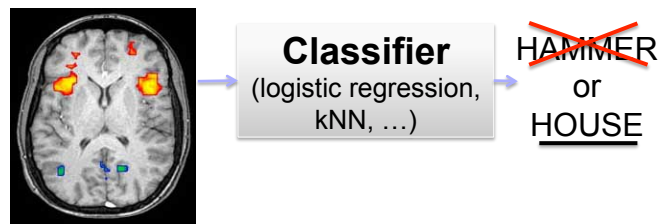
February 14th, 2013

©Emily Fox 2013

12

fMRI Prediction Task

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



©Emily Fox 2013

13

fMRI



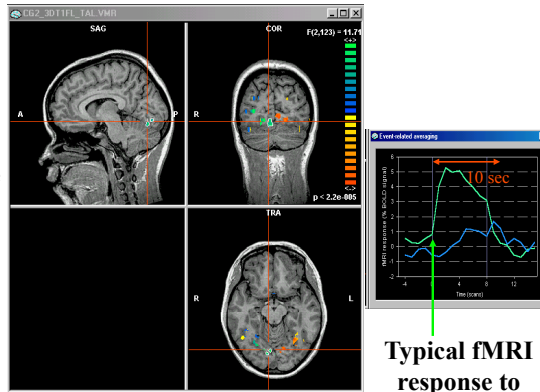
©Emily Fox 2013

14

fMRI

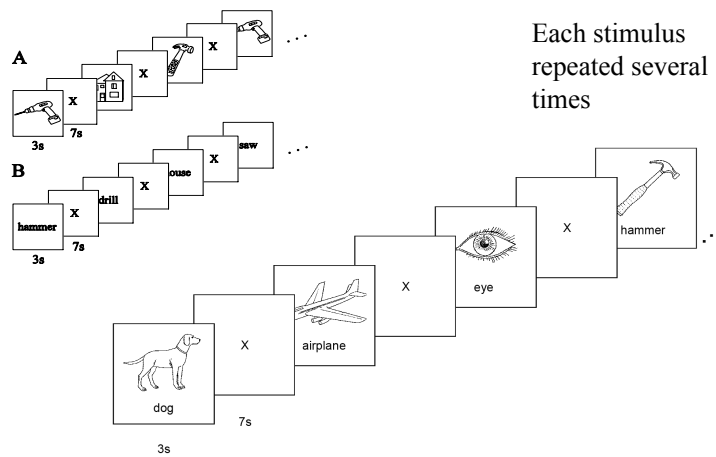
~1 mm resolution
~1 image per sec.
20,000 voxels/image
safe, non-invasive

measures Blood
Oxygen Level
Dependent (BOLD)
response



Typical fMRI
response to
impulse of
neural activity

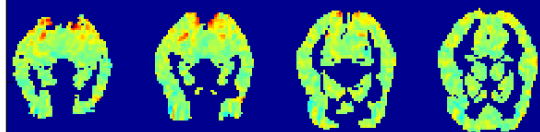
Typical Stimuli



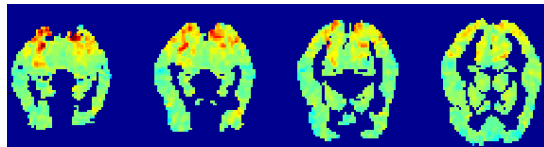
Each stimulus
repeated several
times

fMRI Activation

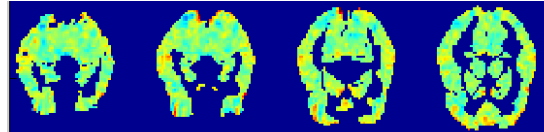
fMRI activation for "bottle":



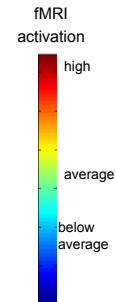
Mean activation averaged over 60 different stimuli:



"bottle" minus mean activation:



bottle

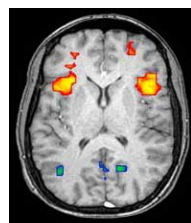


©Emily Fox 2013

17

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



Classifier
(logistic regression,
kNN, ...)

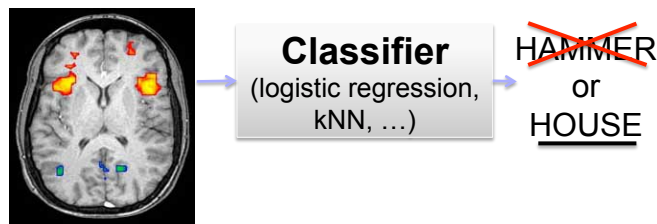
~~HAMMER~~
or
HOUSE

©Emily Fox 2013

18

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

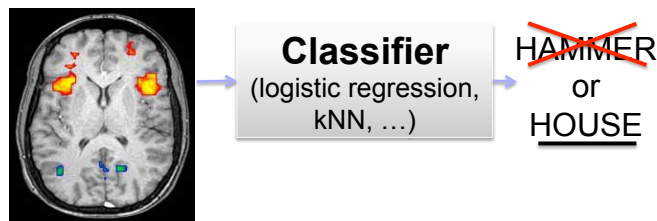


©Emily Fox 2013

19

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?



©Emily Fox 2013

20

Semantic Features

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

Semantic feature values: "airplane"

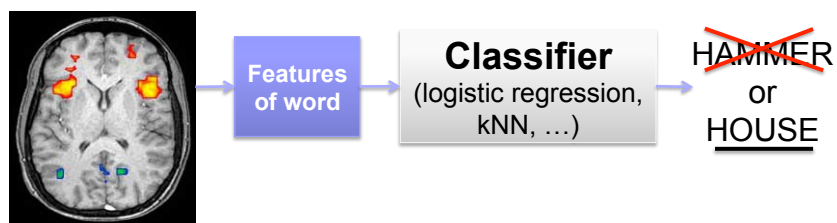
0.8673, ride
0.2891, see
0.2851, say
0.1689, near
0.1228, open
0.0883, hear
0.0771, run
0.0749, lift
...
...
0.0049, smell
0.0010, wear
0.0000, taste
0.0000, rub
0.0000, manipulate

©Emily Fox 2013

21

Zero-Shot Classification

- From training data, learn two mappings:
 - S: input image → semantic features
 - L: semantic features → word
- Can use "cheap" co-occurrence data to help learn L

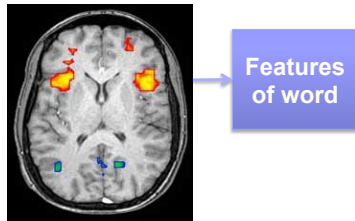


©Emily Fox 2013

22

fMRI Prediction Subtask

- **Goal:** Predict semantic features from fMRI image



©Emily Fox 2013

23

Linear Regression – *review*

- Model:

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

- Minimizing RSS= least squares regression

©Emily Fox 2013

24

Linear Regression – *review*

- Taking the gradient
 - Reformulate objective

- Set gradient = 0

©Emily Fox 2013

25

Ridge Regression

- Ameliorating issues with overfitting:
- New objective:

- Reformulate:

- Set gradient = 0

©Emily Fox 2013

26

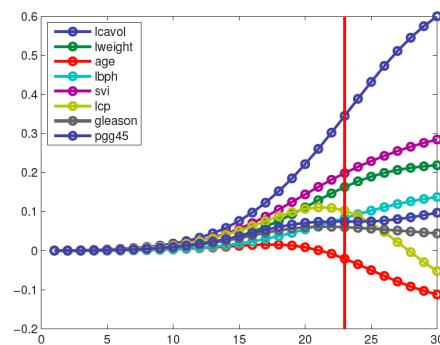
Ridge Regression

- Solution is indexed by the regularization parameter λ
- Larger λ
- Smaller λ
- As $\lambda \rightarrow 0$
- As $\lambda \rightarrow \infty$

©Emily Fox 2013

27

Ridge Coefficient Path



From
Kevin Murphy
textbook

- Typical approach: select λ using cross validation

©Emily Fox 2013

28

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

©Emily Fox 2013

30

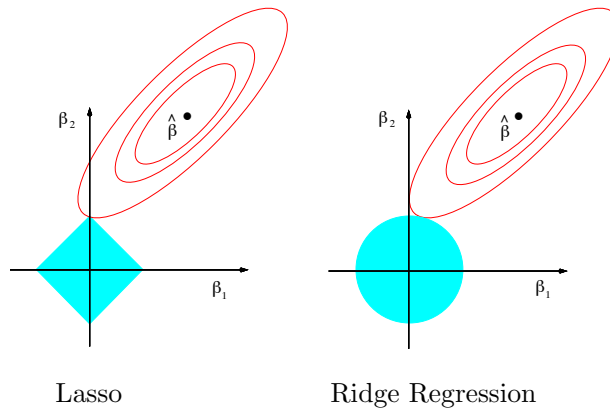
LASSO Regression

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

©Emily Fox 2013

31

Geometric Intuition for Sparsity



©Emily Fox 2013

32

Soft Thresholding

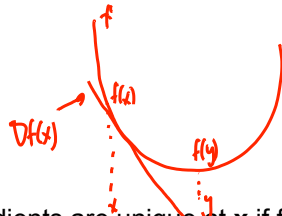
- To see why LASSO results in sparse solutions, look at conditions that must hold at optimum
- L1 penalty $\|\beta\|_1$ is not differentiable whenever $\beta_j = 0$
- Look at subgradient...

©Emily Fox 2013

33

Subgradients of Convex Functions

- Gradients lower bound convex functions:

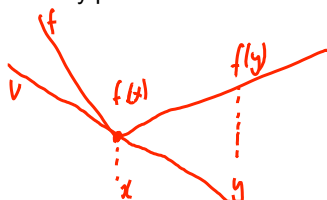


$$f(y) \geq f(x) + \nabla f(x) (y-x)$$

- Gradients are unique at x if function differentiable at x

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



$v \in \partial f(x)$ subgradient
if
 $f(y) \geq f(x) + v \cdot (y-x)$

©Carlos Guestrin 2013

34

Soft Thresholding

- Gradient of RSS term:

- Subgradient of full objective:

©Emily Fox 2013

35

Soft Thresholding

- Set subgradient = 0:

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

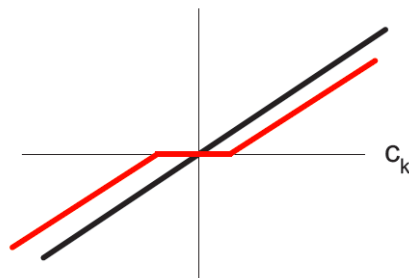
- The value of $c_j = 2 \sum_{i=1}^N x_j^i (y^i - \beta'_{-j} x_{-j}^i)$ constrains β_j

©Emily Fox 2013

36

Soft Thresholding

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

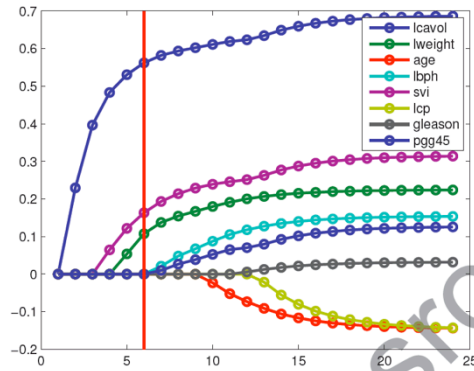


From
Kevin Murphy
textbook

©Emily Fox 2013

37

LASSO Coefficient Path



From Kevin Murphy textbook

©Emily Fox 2013

39

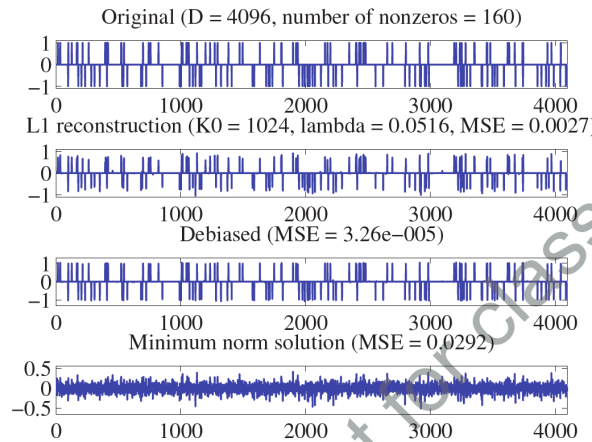
LASSO Example

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	

©Emily Fox 2013

40

Debiasing



From Kevin Murphy textbook

©Emily Fox 2013

41

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 algorithms next time...

©Emily Fox 2013

42

Acknowledgements

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