

Course Overview – Nonparametric Regression and Classification

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 2nd, 2013

©Emily Fox 2013

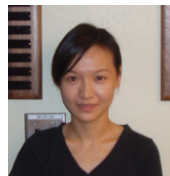
1

Course Staff

- Instructor: **Emily Fox**



- TA: **Shirley You Ren**



©Emily Fox 2013

2

Content: What is the course about?

Course Structure

- 3 Primary Tasks:
 - Regression
 - Classification
 - Density Estimation

- 5 Modules:
 - Nonparametric Preliminaries
 - Splines and Kernels
 - Bayesian Nonparametrics
 - Nonparametrics for Multivariate Covariates
 - Classification

Task 1: Regression

- Assume a sample
- Model:

- Task involves estimating the function f

- Goals of nonparametric approach:
 - Make few assumptions about f
 - Use a large number of parameters, but constrained in some way to avoid overfitting the data
 - Complexity can grow with the sample size

©Emily Fox 2013

5

Task 2: Classification

- Assume a sample $(x_1, Y_1), \dots, (x_n, Y_n)$

- Task involves estimating a predictive model of Y given x

- Goals of nonparametric are as before, but now for link between x and Y with Y discrete-valued

©Emily Fox 2013

6

Task 3: Density Estimation

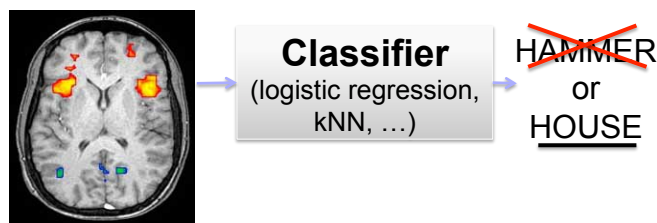
- Assume a sample
- Task involves estimating the density p
- Goals of nonparametric approach are as before, but applied to the estimation of p

©Emily Fox 2013

7

fMRI Prediction Task

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



©Emily Fox 2013

8

fMRI



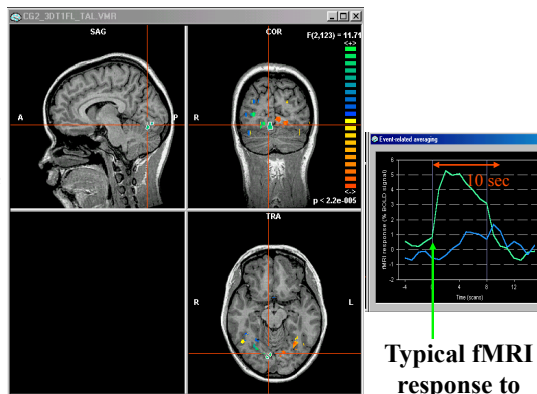
©Emily Fox 2013

9

fMRI

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

measures Blood
Oxygen Level
Dependent (BOLD)
response

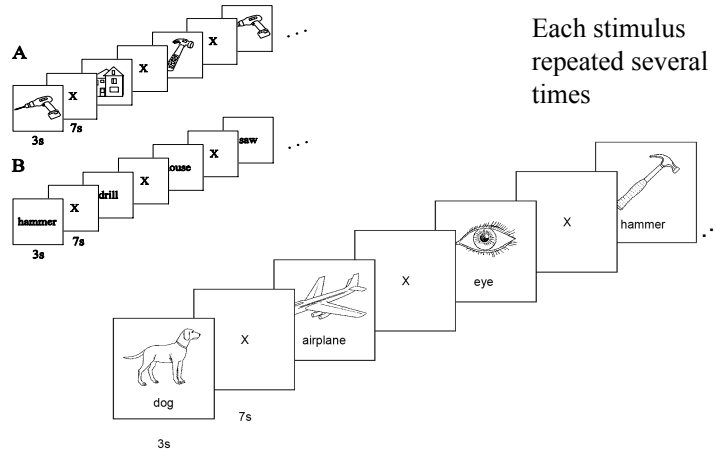


Typical fMRI
response to
impulse of
neural activity

©Emily Fox 2013

10

Typical Stimuli

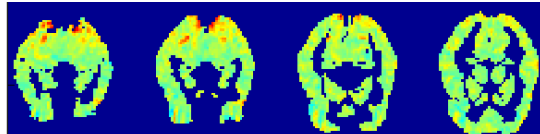


©Emily Fox 2013

11

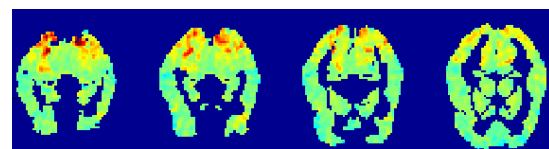
fMRI Activation

fMRI activation for "bottle":

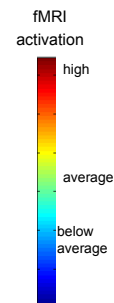
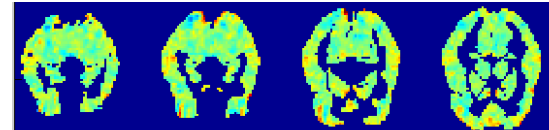


bottle

Mean activation averaged over 60 different stimuli:



"bottle" minus mean activation:

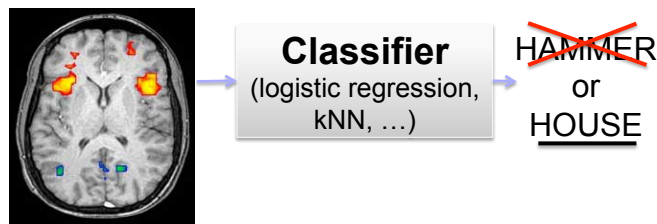


©Emily Fox 2013

12

fMRI Prediction Task

- **Goal:** Predict word stimulus from fMRI image
- **Challenges:**
 - $p \gg n$ (covariate dimension \gg sample size)
 - Cost of fMRI recordings is high
 - Only have a few training examples for each word

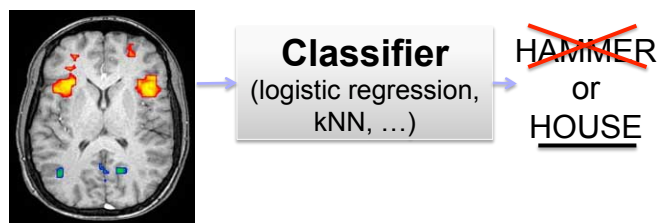


©Emily Fox 2013

13

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

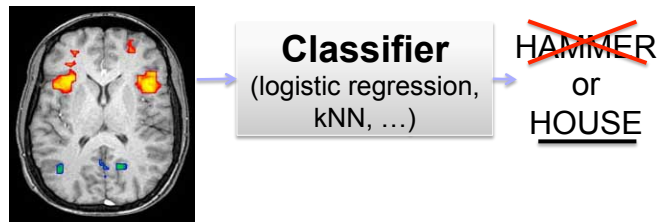


©Emily Fox 2013

14

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

15

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"

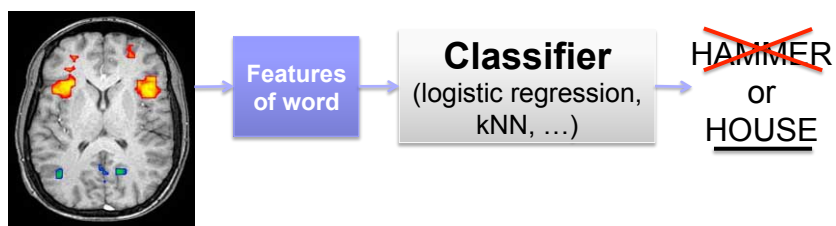
0.8673, ride
0.2891, see
0.2851, say
0.1689, near
0.1228, open
0.0883, hear
0.0771, run
0.0749, lift
...
...
0.0049, smell
0.0010, wear
0.0000, taste
0.0000, rub
0.0000, manipulate

©Emily Fox 2013

16

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



©Emily Fox 2013

17

Assumed Background

- **[Stat 502 and Stat 504] or [Biostat 514 and Biostat 515]**
- **Comfortable with:**
 - Linear algebra
 - Probability
 - R (or Matlab, Python, etc.)
- **Computational and mathematical maturity**

©Emily Fox 2013

18

Logistics: How is the course going to run?

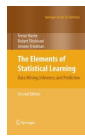
Website and Discussion Board

- Course website:
<http://stat.washington.edu/courses/stat527/s13>

- Catalyst:
 - Used for all discussions
 - Post all questions there (unless personal)
 - See website for sign-up details

Reading

- No required textbook
- Three suggested textbooks (on website):
 - Wakefield, “Bayesian and Frequentist Regression Methods”, Springer 2012
 - Wasserman, “All of Nonparametric Statistics”, Springer 2005
 - Hastie, Tibshirani, Friedman “The Elements of Statistical Learning”, Springer 2009
- Papers linked on course website



Homework

- 7 HWs total
- Assigned and due weekly on *Thursdays*
- Collaboration allowed, but write-ups and coding must be done individually
- Submitted at beginning of class
- Allowed 2 “late days” for entire quarter

Project

- Options:

- Choose project from specified list
- Re-implement existing paper from specified list
- Propose own project idea

- Individual

- New work, but can be connected to research

- Schedule:

- Proposal (1 page) – April 25
- Progress report (3 pages) – May 16
- Poster presentation – June 6
- Final report (8 pages, NIPS format) – June 11

Grading

- HWs 1, 2, 4, 5, 6 (10% each)
- HWs 3, 7 (5%) – short, due dates coincide with project due dates
- Final project (40%)

Support/Resources

■ Office Hours

- TA: W 2-4pm in Padelford A-316
- Emily: Th 11am-12pm in Padelford B-305

■ Recitations

- Optional tutorial/example-based sections will be held *every other* week
- Choose from:
 - Monday, 2-3pm
 - Monday, 5-6pm
 - Tuesday, 4-5pm
- Location TBD

©Emily Fox 2013

25

Module 1: Nonparametric Preliminaries

Intro,
What to Report?

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 2nd, 2013

©Emily Fox 2013

26

The Optimal Prediction

- Assume we *know* the data-generating mechanism
- If our task is prediction, which summary of the distribution $Y | x$ should we report?

- Taking a decision-theoretic framework, consider the ***expected loss***

©Emily Fox 2013

27

Continuous Responses

- Expected loss $E_X \{ E_{Y|X} [L(Y, f(x)) | X = x] \}$

- Example: L_2

Solution:

- Example: L_1

Solution:

- More generally: L_p

©Emily Fox 2013

28

General Responses

- Expected loss $E_X \{ E_{Y|X} [L(Y, f(x)) | X = x] \}$
- Example: log-likelihood

When Gaussian:

When Laplacian:

©Emily Fox 2013

29

Incorporating Models into Prediction

- We don't actually know the data-generating mechanism
- Need an estimator $\hat{f}_n(\cdot)$ based on a random sample Y_1, \dots, Y_n , also known as **training data**

- Statistical models can be used to encode knowledge about aspects of the data-generating mechanism

- Models can provide simplifying assumptions
 - Can help cope with estimation issues due to limited data

©Emily Fox 2013

30

Incorporating Models into Prediction

- Assume some form for how the data are generated
 - E.g., $Y = f(x) + \epsilon$ $E[\epsilon] = 0$ $\text{var}(\epsilon) = \sigma^2$
 - For non-constant variance, can consider GLMs
- Then, typically assume some form for $f(x)$
- Model + loss function \rightarrow some estimator

©Emily Fox 2013

31

Parametric Regression

- *Parametric* inference assumes parametric form for $f(x)$
- Advantages:
 - Efficient estimation
 - Concise summarization
- What is the right parametric form for $f(x)$?

©Emily Fox 2013

32

Goals of Nonparam Regression

- Goals of *nonparametric* inference:
 - Assume little prior knowledge of data-generating mechanism
 - More flexibly model f (i.e., relationship between x and Y)
 - Maintain “reasonable” efficiency of estimation
- Often actually assume parametric forms with large numbers of parameters
 - Constrained to avoid overfitting the data
- Particularly useful when task is prediction
 - Focus on accuracy of prediction rather than parameter values
- Let’s discuss this idea of “complexity” more...

©Emily Fox 2013

33

Model Complexity

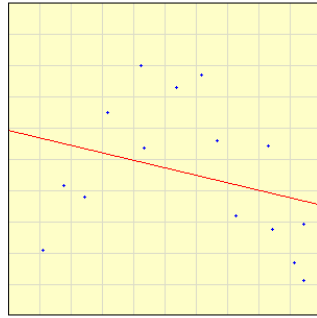
- How complex of a function should we choose?
 - To increase flexibility, using many parameters is attractive
 - However, wide prediction intervals...
 - Leads to wild predictions

©Emily Fox 2013

34

Example: Polynomial Regression

- For added flexibility, allow for high order polynomial, right?



Select points by clicking on the graph or press [Example](#)

Degree of polynomial: Fit Y to X
 Fit X to Y

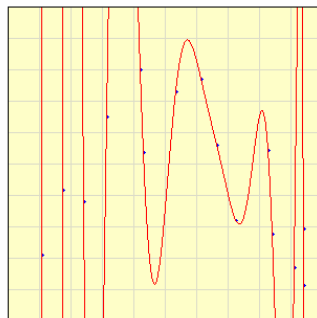
[Calculate](#) [View Polynomial](#) [Reset](#)

©Emily Fox 2013

35

Example: Polynomial Regression

- For added flexibility, allow for high order polynomial, right?



Select points by clicking on the graph or press [Example](#)

Degree of polynomial: Fit Y to X
 Fit X to Y

[Calculate](#) [View Polynomial](#) [Reset](#)

©Emily Fox 2013

36

Measuring Predictive Performance

- Assume estimate $\hat{f}_n(\cdot)$ based on training data y_1, \dots, y_n
- The **generalization error** provides a measure of predictive performance

$$GE(\hat{f}_n) = E_{Y,X} [L(Y, \hat{f}_n(X))]$$

©Emily Fox 2013

37

Measuring Predictive Performance

- Assume L_2 loss
- Averaging over repeat training sets $\mathbf{Y}_n = Y_1, \dots, Y_n$ we get the **predictive risk** at x^*

$$E_{Y^*, \mathbf{Y}_n} [(Y^* - \hat{f}_n(x^*))^2] =$$

- Recall $MSE[\hat{f}_n(x)] = \text{bias}(\hat{f}_n(x))^2 + \text{var}(\hat{f}_n(x))$

©Emily Fox 2013

38

Measuring Predictive Performance

- Finally, let's average over covariates x
 - *Integrated MSE*
 - *Average MSE*
- Note: ***avg. pred. risk*** = $\sigma^2 + \text{avg. MSE}$

©Emily Fox 2013

39

Bias-Variance Tradeoff

- Minimizing risk = balancing bias and variance

- Note: *f(x)* is unknown, so cannot actually compute MSE

©Emily Fox 2013

40

More on Nonparam Regression

- Often framed as learning functions with a complexity penalty
 - Regular behavior in small neighborhoods of the input
 - E.g., locally linear or low-order polynomial...estimator results from averaging over these local fits
- Choice of neighborhood = strength of constraint
 - Large neighborhood can lead to linear fit (very restrictive) whereas small neighborhoods can lead to interpolation (no restriction)

©Emily Fox 2013

41

More on Nonparam Regression

- Different restrictions lead to different nonparametric approaches
 - Roughness penalty → *splines*
 - Weighting data locally → *kernel methods*
 - Etc.
- Each method has associated *smoothing* or *complexity* param
 - Magnitude of penalty
 - Width of kernel (defining “local”)
 - Number of basis functions
 - ...
- Bias-variance tradeoff
- Will explore methods for choosing smoothing parameters

©Emily Fox 2013

42

Module 1: Nonparametric Preliminaries

Review of Regression, Linear Smoothers

STAT/BIOSTAT 527, University of Washington

Emily Fox

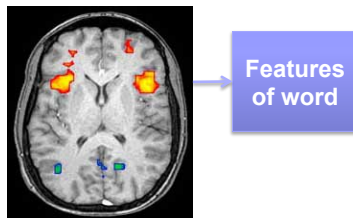
April 2nd, 2013

©Emily Fox 2013

43

fMRI Prediction Subtask

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



©Emily Fox 2013

44

Linear Regression – *review*

- Model:
- *Design matrix*:
- Rewrite in matrix form:

©Emily Fox 2013

45

Linear Regression – *review*

- Least squares estimation:
 - Minimize *residual sum of squares*
 - Take gradient and set = 0
- In Gaussian case, LS est. = maximum likelihood est.

©Emily Fox 2013

46

Linear Regression – *review*

- **Fitted values**
- Number of parameters
- For any x , we can write

©Emily Fox 2013

47

Linear Smoothers

- Definition:
 \hat{f}_n of f is a **linear smoother** if, for each x , there exists
$$\ell(x) = (\ell_1(x), \dots, \ell_n(x))^T$$
such that
- Matrix form
 - Fitted values
 - Smoothing or “hat” matrix
- Effective degrees of freedom:

©Emily Fox 2013

48

Linear Smoothers

- Note 1:

A linear smoother does **not** imply that $f(x)$ is linear in x

- Note 2:

If $Y_i = c$ for all i , then $\hat{f}_n(x) = c$ for all x