# Course Overview – Nonparametric Regression and Classification

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 2nd, 2013

---

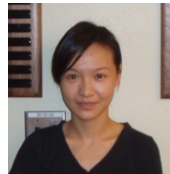# Course Staff

- Instructor: **Emily Fox**



- TA: **Shirley You Ren**

# 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 $(x_1, Y_1), \ldots, (x_n, Y_n)$
- Model: $Y_i = f(x_i) + \epsilon_i \qquad E[\epsilon_i] = 0$

  $\uparrow$ unknown



- Task involves estimating the function $f$

  What $f$ should I use?

- 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

5

---

# Task 2: Classification

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

  $Y_i \in \{1, \ldots, k\}$

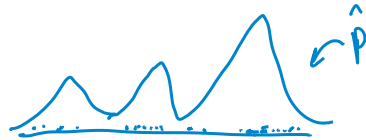  $\uparrow$ # classes

  2 classes
  +
  o



- 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

6

# Task 3: Density Estimation

- Assume a sample $X_1, \ldots, X_n \sim p$

  *unknown*

  $\hat{p}$

- Task involves estimating the density $p$

- Goals of nonparametric approach are as before, but applied to the estimation of $p$
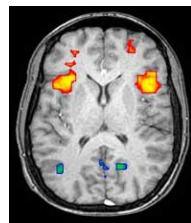
©Emily Fox 2013

7

# fMRI Prediction Task

*cool task involving both reg. + class.*

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

*Can we read your brain?*



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

HAMMER
or
HOUSE

©Emily Fox 2013

8

4

# fMRI

9

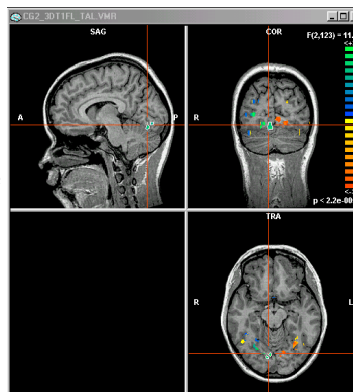# fMRI

*very high*
*pretty slow*

**~1 mm resolution**
**~1 image per sec.**

**20,000 voxels/image**

**safe, non-invasive**
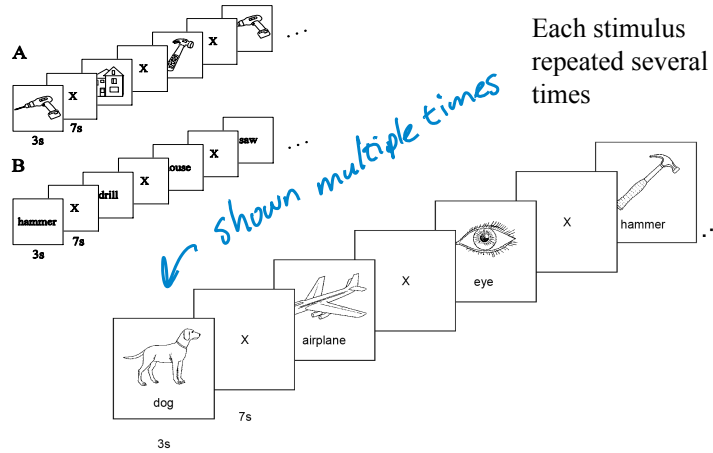
**measures Blood Oxygen Level Dependent (BOLD) response**

**Typical fMRI response to impulse of neural activity**

10

5

# Typical Stimuli



Each stimulus repeated several times

*← shown multiple times*

11

# fMRI Activation

fMRI activation for "bottle":



*← stimulus*

bottle

Mean activation averaged over 60 different stimuli:



fMRI activation

high

average

below average

"bottle" minus mean activation:

*Is this enough?*

12

# fMRI Prediction Task

- **Goal:** Predict word stimulus from fMRI image
- **Challenges:** ✓ — # of voxels
    - ☐ p >> n (covariate dimension >> sample size)
    - ☐ Cost of fMRI recordings is high
    - ☐ Only have a few training examples for each word

**Classifier**
(logistic regression, kNN, …) → ~~HAMMER~~ or <u>HOUSE</u>

©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

*Never showed the word "giraffe"*

**Classifier**
(logistic regression, kNN, …) → ~~HAMMER~~ or <u>HOUSE</u>

©Emily Fox 2013                    14

# Zero-Shot Classification

*(handwritten: many docs containing "giraffe" also contain "neck" "zoo" ...)*
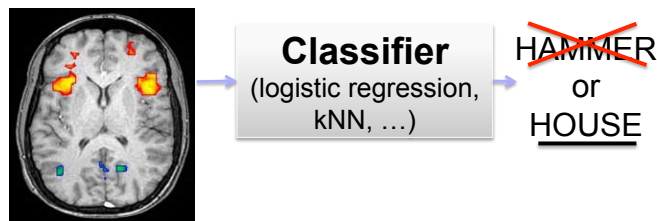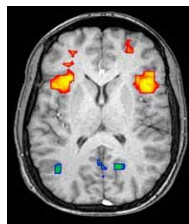
- **Goal:** Classify words not in the training set
- **Challenges:**
  - Cost of fMRI recordings is high
  - Can't get recordings for every word in the vocabulary
- We don't have many brain images, but we have a lot of info about the words and how they relate (co-occurrence, etc.)
- How do we utilize this "cheap" information?



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

~~HAMMER~~
or
HOUSE

15

---

# Semantic Features

*(handwritten: Google Trillion word Corpus)*

| 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 |

16

# Zero-Shot Classification

- From training data, learn two mappings: $A = \{ \blacksquare \rightarrow \text{"dog"} \}$ *few*
  - S: input image → semantic features
  - L: semantic features → word

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

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

  *Training* $\{ \blacksquare \rightarrow [\vdots] \rightarrow \text{"dog"} \}$ *uses* $A + B$ — *n examples, n small*



**Features of word** → **Classifier** (logistic regression, kNN, …) → ~~HAMMER~~ or <u>HOUSE</u>

*"giraffe"*

*Predict* $\blacksquare \xrightarrow{S} [\vdots] \xrightarrow{L} \text{"giraffe"}$

*using training data*

*Questions: -S? -L?*

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**

18

# Logistics: How is the course going to run?

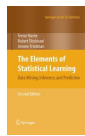# Website and Discussion Board

- 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

23

# Grading

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

24

# 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

25

---

# Module 1: Nonparametric Preliminaries

## Intro,
## What to Report?

STAT/BIOSTAT 527, University of Washington

Emily Fox

April 2nd, 2013

26

# The Optimal Prediction

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

  For $x$, what fcn $f(x)$ should we choose to predict $Y$ if we can choose any $f(\cdot)$

- Taking a decision-theoretic framework, consider the **expected loss**  predictions are penalized by $L(\cdot, \cdot)$

$$E_{X,Y}\left[L(Y, f(X))\right] = E_X\left\{E_{Y|x}\left[L(Y, f(x)) \mid X = x\right]\right\}$$

- $\hat{f}(\cdot)$ should min
- can min. pointwise

27

---

# Continuous Responses

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

- Example: $L_2$  $L(Y, f(x)) = (Y - f(x))^2$

  Solution: $\hat{f}(x) = E[Y|x]$

  Proofs: HW

- Example: $L_1$  $L(Y, f(x)) = |Y - f(x)|$

  Solution: $\hat{f}(x) = \text{median}(Y|x)$

- More generally: $L_p$  $L(Y, f(x)) = \left\{\int |Y - f(x)|^p\right\}^{1/p}$

©Emily Fox 2013

28

14

# General Responses

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

- Example: log-likelihood $\quad L(Y, f(x)) = -2 \log p(Y \mid f(x))$

  When Gaussian: $Y|x \sim N(f(x), \sigma^2)$

  $$\rightarrow L(Y, f(x)) = \log(2\pi\sigma^2) + (Y - f(x))^2 / \sigma^2$$

  $$\rightarrow \hat{f}(x) = E[Y|x]$$

  When Laplacian: $Y|x \sim Lap[f(x), \phi]$

  $$\rightarrow \hat{f}(x) = median(Y|x)$$

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 $Y_1, \ldots, Y_n$, also known as **training data**

  e.g. Est. $E[Y|x]$, but how? Typically don't have mult. obs. at a given x. Maybe $\hat{f}_n(x) = Avg(Y_i \mid X_i \in Nbhd(x))$

  Can be problematic if not many obs.

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

  e.g. Assume linear form, then we know how to approx $E[Y|x]$ s.t. a linear const. on f

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

30

# Incorporating Models into Prediction

- Assume some form for how the data are generated
  - E.g., $Y = f(x) + \epsilon \qquad E[\epsilon] = 0 \quad \text{var}(\epsilon) = \sigma^2$

  $$\Rightarrow f(x) = E[Y|x]$$

  *how est. $f$ ?*

  - For non-constant variance, can consider GLMs
- Then, typically assume some form for *f(x)*

  $$e.g., \quad f(x) = \beta^T x$$

- Model + loss function → some estimator  *r.v.*

  $$e.g. \ L_2 \ \rightarrow \ \hat{\beta} = \left( E(XX^T) \right)^{-1} E(XY)$$

  $$approx. \quad \hat{\beta}_n = (X^T X)^{-1} X^T y \quad \leftarrow obs.\ vec.$$

  *matrix of covariates*

# Parametric Regression

- *Parametric* inference assumes parametric form for $f(x)$

  $$e.g. \quad f(x) = \beta^T X$$

  *$f(\cdot)$ is indexed by param $\beta$*

- Advantages:
  - Efficient estimation  ← *e.g. LS est. of $\beta$, $\hat{\beta}_n$,*
  - Concise summarization    *leads to an est. $\hat{f}_n$ of $f$*

- What is the right parametric form for $f(x)$?

  *Should it change w/ sample size?*

16

# 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…
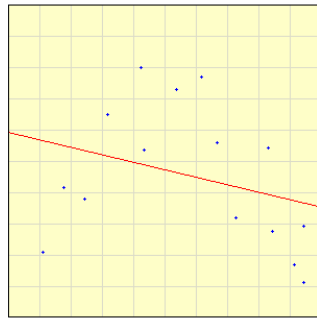
33

# Model Complexity

- How complex of a function should we choose?

  - To increase flexibility, using many parameters is attractive

    Reduce bias

  - However, wide prediction intervals…

    Fixed dataset contains a limited amt of info

  - Leads to wild predictions

34

17

# Example: Polynomial Regression

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

$$y_i = \sum_{j=0}^{p} \beta_j X_i^j + \epsilon_i$$

Not always good to add params

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

Degree of polynomial: [1 ▾]  ● Fit Y to X   ○ Fit X to Y
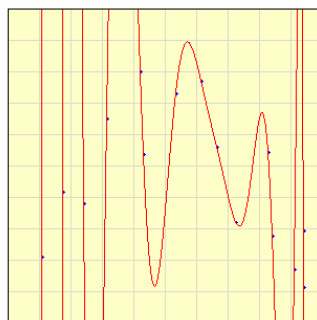
[Calculate] [View Polynomial] [Reset]

35

---

# Example: Polynomial Regression

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

sensitive to small changes in data

High order = low bias, but high var

How do we assess an estimator fn ?

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

Degree of polynomial: [13 ▾]  ● Fit Y to X   ○ Fit X to Y

[Calculate] [View Polynomial] [Reset]

36

18

# Measuring Predictive Performance

- Assume estimate $\hat{f}_n(\cdot)$ based on training data $y_1,..., y_n$

  $\uparrow$ fixed

- The **generalization error** provides a measure of predictive performance

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

  $\uparrow$ fixed

# Measuring Predictive Performance

- Assume $L_2$ loss    $Y = f(x) + \epsilon$  $\bigstar$  $E[\epsilon] = 0$  $var(\epsilon) = \delta^2$
- Averaging over repeat training sets $\mathbf{Y}_n = Y_1,..., Y_n$ we get the **predictive risk** at $x^*$

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

test  training  $\underset{\text{training data } \mathbf{Y}_n}{\underbrace{\text{fcn of}}}$

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

$$\underbrace{\qquad}_{MSE(\hat{f}_n^{(x)})} \qquad O$$

$$= \delta^2 + MSE\left(\hat{f}_n(x^*)\right) \qquad \checkmark$$

$\uparrow$ "irreducible error"  "risk"

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

# Measuring Predictive Performance

- Finally, let's average over covariates $x$

  - **Integrated MSE**  $\int MSE(\hat{f}_n(x)) p(x)dx$

    *summary over all inputs*

  - **Average MSE**  $\frac{1}{n}\sum_{i=1}^{n} MSE(\hat{f}_n(x_i))$

    *Monte Carlo est.*
    $x_i \sim P$

- Note:  **avg. pred. risk =** $\sigma^2$ **+ avg. MSE**

$$\left[\frac{1}{n}\sum_{i=1}^{n} E_{Y_n, Y_n^*}\left[(Y_i^* - \hat{f}(x_i))^2\right]\right]$$

*training*    *new obs.* $Y_n^* = Y_1^*, \ldots, Y_n^*$

39

# Bias-Variance Tradeoff

*recall polynomial reg. example*

- Minimizing risk = balancing bias and variance



*MSE*  *Var*  *bias*

*MSE is a fcn of f(·)*    *inc.* *model complexity* →

*optimal solution*

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

40

20

# 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)

---

# 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