

Module 1: Nonparametric Preliminaries

Review of Regression, Linear Smoothers

STAT/BIOSTAT 527, University of Washington

Emily Fox

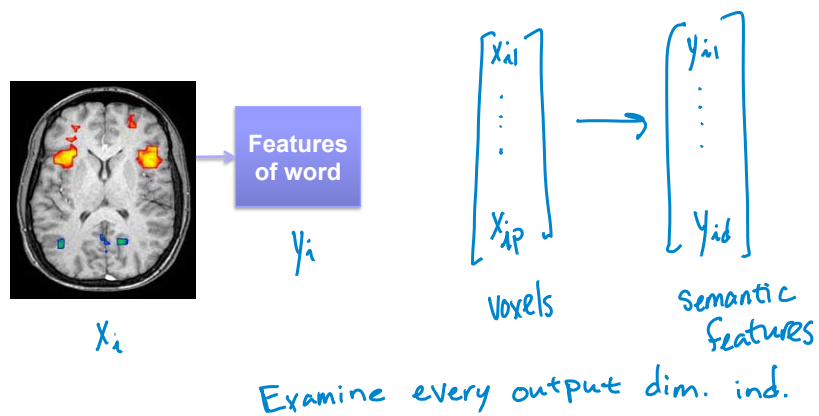
April 4th, 2013

©Emily Fox 2013

1

fMRI Prediction Subtask

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



©Emily Fox 2013

2

Linear Regression – review

- Model:
$$y_i = \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i \quad i=1, \dots, n$$

$$= X_i^T \beta + \epsilon_i$$

$$E[\epsilon_i] = 0 \quad \text{var}(\epsilon_i) = \sigma^2$$

$$X_{i1} = 1 \quad \text{for intercept}$$
- Design matrix:
$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

$$\beta = (\beta_1, \dots, \beta_p)^T$$
- Rewrite in matrix form:
$$Y = X\beta + \epsilon$$

©Emily Fox 2013

3

Linear Regression – review

- Least squares estimation:
 - Minimize **residual sum of squares**

$$\hat{\beta} = \underset{\beta}{\text{argmin}} (Y - X\beta)^T (Y - X\beta) = \sum_{i=1}^n (y_i - \sum_j x_{ij} \beta_j)^2$$

$$\frac{1}{2} \text{RSS}(\beta) = \frac{1}{2} \beta^T (X^T X) \beta - \beta^T X^T Y + \text{const.}$$

$$\text{RSS}(\beta)$$
 - Take gradient and set = 0

$$\nabla_{\beta} \frac{1}{2} \text{RSS}(\beta) = X^T X \beta - X^T Y = 0$$

$$\Rightarrow \hat{\beta}^{\text{LS}} = (X^T X)^{-1} X^T Y$$
- In Gaussian case, LS est. = maximum likelihood est.

©Emily Fox 2013

4

Linear Regression – review $f(x) = X^T \beta$

- **Fitted values**

$$\hat{f}_n = X \hat{\beta}^{LS} = LY$$

$$L = X(X^T X)^{-1} X^T$$

"hat matrix"

- Number of parameters

$$p = \text{tr}(L)$$

proof: $\text{tr}(X(X^T X)^{-1} X^T)$
 $= \text{tr}((X^T X)^{-1} X^T X) = \text{tr}(I_p)$

- For any x , we can write

$$\hat{f}_n(x) = l(x)^T y = \sum_{i=1}^n l_i(x) y_i$$

where $l(x) = x(X^T X)^{-1} X^T$
 ↑ all x_{ij} 's from training data

©Emily Fox 2013

5

Linear Smoothers

- **Definition:**

\hat{f}_n of f is a **linear smoother** if, for each x , there exists

$$l(x) = (l_1(x), \dots, l_n(x))^T \quad \text{with} \quad \sum_{i=1}^n l_i(x) = 1$$

such that

$$\hat{f}_n(x) = \sum_{i=1}^n l_i(x) y_i$$

- **Matrix form**

- Fitted values

$$\hat{f} = LY$$

- Smoothing or "hat" matrix

$$L_{ij} = l_j(x)$$

- Effective degrees of freedom:

$$v = \text{tr}(L)$$

©Emily Fox 2013

6

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

bc $\sum_{i=1}^n \hat{l}_i(x) = 1$



©Emily Fox 2013

7

Module 1: Nonparametric Preliminaries

Overfitting,
Ridge Regression,
LASSO

STAT/BIOSTAT 527, University of Washington

Emily Fox

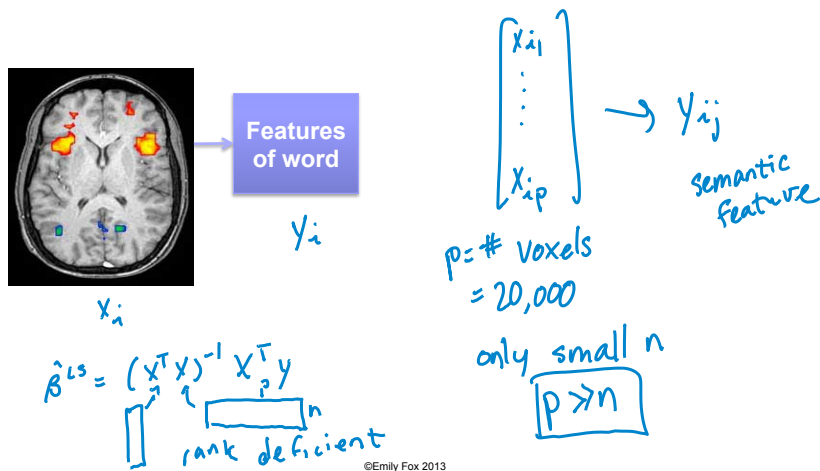
April 4th, 2013

©Emily Fox 2013

8

fMRI Prediction Subtask

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



©Emily Fox 2013

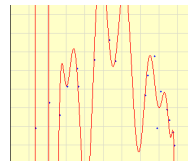
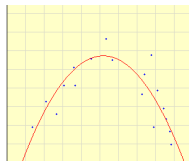
9

Regularization in Linear Regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$

$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



even for $n \gg p$, p large

- **Regularized** or **penalized** regression aims to impose a “complexity” penalty by penalizing large weights
 - “Shrinkage” method

©Carlos Guestrin 2005-2009

10

Ridge Regression

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

- New objective:

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

↑ don't penalize the intercept
 ↑ $\beta^T \beta$
 ↑ strength of the penalty

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

©Emily Fox 2013

11

Ridge Regression

- New objective:

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \sum_{i=1}^n (y_i - (\beta_0 + \beta^T x_i))^2 + \lambda \|\beta\|_2^2$$

$f(\beta)$

- Reformulate:

$$\frac{1}{2} F(\beta) = \frac{1}{2} \beta^T (X^T X) \beta - \beta^T X^T y + \text{const.} + \frac{1}{2} \lambda \beta^T \beta$$

RSS(β) as before

- Set gradient = 0

$$= \frac{1}{2} \beta^T (X^T X + \lambda I) \beta - \beta^T X^T y + \text{const.}$$

as before

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

- Linear smoother!!

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

©Emily Fox 2013

12

Ridge Regression

- Solution is indexed by the regularization parameter λ
- Larger λ *high reg.*
- Smaller λ *low reg.*
- As $\lambda \rightarrow 0$ $\hat{\beta}^{ridge} \rightarrow \hat{\beta}^{LS}$
- As $\lambda \rightarrow \infty$ $\hat{\beta}^{ridge} \rightarrow 0$

©Emily Fox 2013

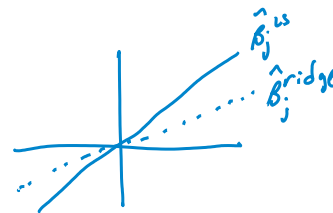
13

Shrinkage Properties

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

- If orthogonal covariates $X^T X = I$

$$\hat{\beta}_j^{ridge} = \frac{\hat{\beta}_j^{LS}}{1+\lambda} \quad v_j$$



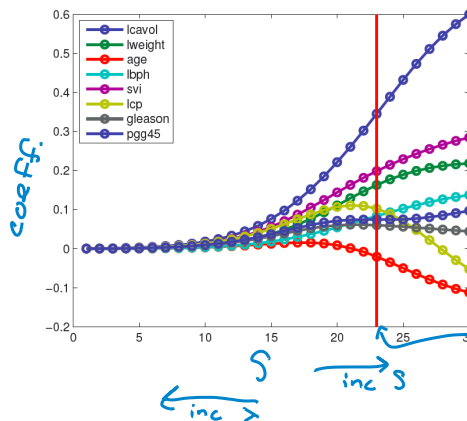
- Effective degrees of freedom:

$$\nu = \text{tr}(L) = \text{tr}(X(X^T X + \lambda I)^{-1} X^T)$$

©Emily Fox 2013

14

Ridge Coefficient Path



From Kevin Murphy textbook

$$\|\beta\|_2^2 = S$$

CV solution

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

A Bayesian Formulation

- Consider a model with likelihood

$$y_i | \beta \sim N(\beta_0 + x_i^T \beta, \sigma^2)$$

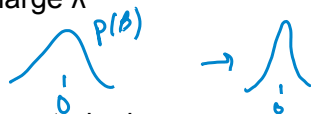
$$\text{if } \epsilon \sim N(0, \sigma^2)$$

and prior

$$\beta \sim N\left(0, \frac{\sigma^2}{\lambda} I_p\right)$$

$$\beta_j \sim N\left(0, \frac{\sigma^2}{\lambda}\right)$$

- For large λ



prior is peaked around $\beta=0$
 \Leftrightarrow penalizing β far from 0

- The posterior is

$$\beta | y \sim N\left(\hat{\beta}^{\text{ridge}}, \sigma^2 (X^T X + \lambda I)^{-1} X^T X \sigma^2 (X^T X + \lambda I)^{-1}\right)$$

easy to show
 $\text{Var}(\hat{\beta}^{\text{ridge}})$

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

discrete

2^p subsets of predictors... can't do this

- Stepwise procedures are sensitive to data perturbations and often include features with negligible improvement in fit

← greedy, 3 backtracking alg.

- not min this obj.
- coeff. sensitive to what's in c. in the model

- Try new penalty: Penalize non-zero weights

□ Penalty: $\|B\|_1 = \sum_j |B_j|$

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

LASSO Regression

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

- New objective:

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

RSS(B)



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

LASSO Solutions

- The LASSO solution is **nonlinear** in y ...*not a linear smoother*
 - Degrees of freedom cannot be computed as before
 - Many recent studies on this (e.g., Zou et al. 2007, Tibshirani & Taylor 2011)
 - Standard errors via the bootstrap

- Efficient algorithms exist for solving
 - Will return to this in a few slides

Geometric Intuition for Sparsity

