

Module 4: Coping with Multiple Predictors

Multidimensional Splines

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 6th, 2014

©Emily Fox 2014

1

Nonparam. Multiple Regression

- We now consider a d -dimensional covariate x_i
$$x_i = (x_{i1}, \dots, x_{id}) \quad i=1, \dots, n$$
- In its most general form, the regression equation then takes the form
$$y = f(x_1, \dots, x_d) + \varepsilon$$

or, for GLMs,

$$g(E(y)) = f(x_1, \dots, x_d)$$
- In principle, all of the methods we have discussed so far carry over to this case rather straightforwardly
- Unfortunately, the risk of the nonparametric estimator increases rapidly with covariate dimension d .
$$\text{pred risk} = \text{MSE} + \sigma^2$$

©Emily Fox 2014

2

Curse of Dimensionality

- To maintain a fixed level of accuracy for a given nonparametric estimator, the sample size must increase exponentially in d

- Set $MSE = \delta$

$$n \propto \left(\frac{c}{\delta}\right)^{d/4}, c > 0$$

- Why? Using data in local nbhd
 - In high dim, few points in any nbhd

- Consider example with n uniformly distributed points in $[-1, 1]^d$

- $d=1$:

$$\text{in } [-0.1, 0.1] \sim n \times \left(\frac{1}{10}\right)$$

- $d=10$

$$\text{in } [-0.1, 0.1]^{10} = \frac{n}{10^{10}} \text{ in interval}$$

$$\sim n \times \left(\frac{1}{10}\right)^{10}$$

$$= \frac{n}{10,000,000,000}$$

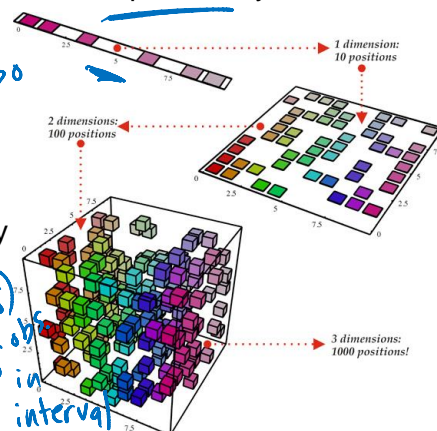


Figure from Yoshua Bengio's website

©Emily Fox 2014

3

Natural Thin Plate Splines

- One-dimensional smoothing splines (obtained via regularization) can be extended to the multivariate setting as the solution to

$$\min_f \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda J(f)$$

\uparrow
 $x_i \in \mathbb{R}^d$

- Recall roughness penalty in 1d

$$J(f) = \int f''(x)^2 dx$$

- The natural 2d extension to penalize rapid variation in either dim is

$$J(f) = \iint_{\mathbb{R}^2} \left[\left(\frac{\partial^2 f(x)}{\partial x_1^2} \right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_2^2} \right)^2 + 2 \left(\frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 \right] dx_1 dx_2$$

- Is the penalty affected by rotation or translation in \mathbb{R}^2 ?

(can be extended for $d > 2$) No

©Emily Fox 2014

4

Natural Thin Plate Splines

$$\min_f \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda J(f)$$

$x_i \in \mathbb{R}^2$

$$J(f) = \int \int_{\mathbb{R}^2} \left[\left(\frac{\partial^2 f(x)}{\partial x_1^2} \right)^2 + 2 \left(\frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_2^2} \right)^2 \right] dx_1 dx_2$$

"bending energy"

- Solution: Unique minimizer is the natural thin plate spline with knots at the x_{ij}
- Proof: See Green and Silverman (1994) and Duchon (1977)
- Similar properties and intuition as in 1d:
 - As $\lambda \rightarrow 0$, sol'n approaches an interpolator
 - As $\lambda \rightarrow \infty$, LS plane (no 2nd derivative)

©Emily Fox 2014

5

Natural Thin Plate Splines

$$\min_f \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda J(f)$$

$$J(f) = \int \int_{\mathbb{R}^2} \left[\left(\frac{\partial^2 f(x)}{\partial x_1^2} \right)^2 + 2 \left(\frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_2^2} \right)^2 \right] dx_1 dx_2$$

- Solution: natural thin plate spline with knots at the x_{ij}
- For general λ , solution is a linear basis expansion of the form

$$f(x) = \beta_0 + \beta^T x + \sum_{j=1}^n b_j h_j(x)$$

with

$$h_j(x) = \|x - x_j\|^2 \log \|x - x_j\|$$

$x \in \mathbb{R}^2$

Radial Basis Function

- Interpretation: We take an elastic flat plate that interpolates points (x_i, y_i) and penalize its "bending energy"

©Emily Fox 2014

6

Natural Thin Plate Splines

$$f(x) = \beta_0 + \beta^T x + \sum_{j=1}^n b_j h_j(x)$$

- Coefficients are found via standard penalized LS

$$\min_{\beta, b} (y - X\beta - Eb)^T (y - X\beta - Eb) + \lambda b^T E b$$

$f(x)$

$$\text{s.t. } \sum_i b_i = \sum_i b_i x_{i1} = \sum_i b_i x_{i2} = 0$$

$\sum b_i = 0$ ensures finite penalty

penalty

$$E_{ij} = \|x_i - x_j\|^2 \log \|x_i - x_j\|$$

- Interpretation: We take an elastic flat plate that interpolates points (x_i, y_i) and penalize its “bending energy”

©Emily Fox 2014

7

Complexity of Thin Plate Splines

- Natural thin plate splines place knots at every location x_{ij}

lots of knots

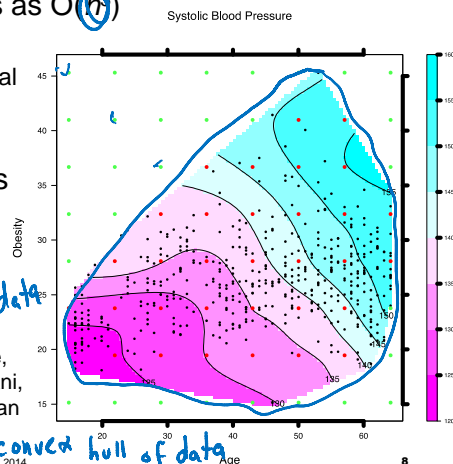
- Computational complexity scales as $O(n^3)$

- Can get away with fewer knots
- If we use K knots, then computational complexity reduces to $O(nK^2 + K^3)$

$K \ll n$

- Can choose some lattice of knots

From Hastie, Tibshirani, Friedman book
= ignore heart disease data outside convex hull of data



©Emily Fox 2014

8

Thin Plate Regression Splines

- Thin plate regression splines truncate the “wiggly” basis b ($K < n$)
- Let $E = UDU^T$ *eigen decomp*
 \uparrow *eigvec* \uparrow *diag matrix of eigenvalues (ordered)*
- Grab out largest k eigenvalues and eigenvectors
 $\downarrow D_k = k \times k$ submatrix of D . $U_k =$ 1st k columns of U .
- Define $b = U_k b_k$
- Minimize

$$\min_{\beta, b_k} (y - X\beta - U_k D_k b_k)^T (y - X\beta - U_k D_k b_k) + \lambda b_k^T D_k b_k$$

$$X^T U_k b_k = 0$$
 $E_k b = (U_k D_k U_k^T) b_k$
- Optimal approximation of thin plate splines using low rank basis
- Retain advantages of (i) no choice of knots, (ii) rotation invariance
- See Wood (2006) for more details

©Emily Fox 2014

9

Tensor Product Splines

- Again, assume x in \mathbb{R}^2 (k if $d=2$)
- Instead of thin plate splines, consider modeling $f(x)$ as follows
- Suppose for each dimension (x_1, x_2) we have a basis of functions

$$h_{1k}(x_1) \quad k=1, \dots, M_1$$

$$h_{2k}(x_2) \quad k=1, \dots, M_2$$

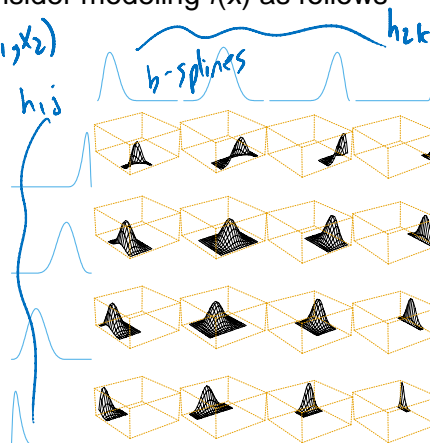
- Then the $M_1 \times M_2$ dimensional **tensor product basis** is

$$g_{jk}(x) = h_{1j}(x_1) h_{2k}(x_2)$$

\uparrow
 \mathbb{R}^2

$$j=1, \dots, M_1$$

$$k=1, \dots, M_2$$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

10

Tensor Product Splines

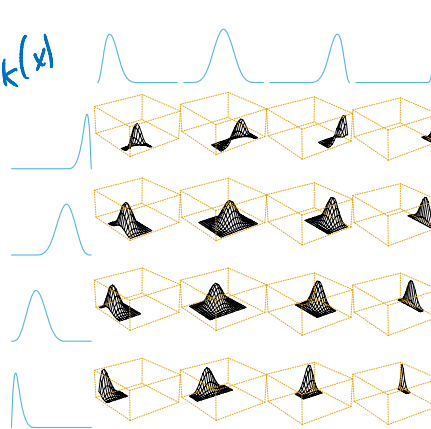
- We use this tensor product basis

$$g_{jk}(x) = h_{1j}(x_1)h_{2k}(x_2)$$

to model $f(x)$

$$f(x) = \sum_{j=1}^{M_1} \sum_{k=1}^{M_2} \theta_{jk} g_{jk}(x)$$

- This formulation extends (in theory) to any dimension d
- Note that as the dimension of the basis grows exponentially with the input dimension d



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

11

Tensor Product Splines Example

- Linear spline basis with L_1 truncated lines for x_1 and L_2 for x_2

$$1, x_1, (x_1 - \xi_{11})_+, \dots, (x_1 - \xi_{1L_1})_+$$

$$1, x_2, (x_2 - \xi_{21})_+, \dots, (x_2 - \xi_{2L_2})_+$$

- Then, the tensor product expansion is

$$f(x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \sum_{l_1=1}^{L_1} b_{l_1} (x_1 - \xi_{1l_1})_+ + \sum_{l_2=1}^{L_2} b_{l_2} (x_2 - \xi_{2l_2})_+ + \sum_{l_1=1}^{L_1} c_{l_1} x_2 (x_1 - \xi_{1l_1})_+ + \sum_{l_2=1}^{L_2} c_{l_2} x_1 (x_2 - \xi_{2l_2})_+$$

- Number of parameters:

$$(L_1 + 2) + (L_2 + 2) + L_1 + L_2 + L_1 + L_2 + L_1 L_2 + \sum_{l_1=1}^{L_1} \sum_{l_2=1}^{L_2} d_{l_1 l_2} (x_1 - \xi_{1l_1})_+ (x_2 - \xi_{2l_2})_+$$

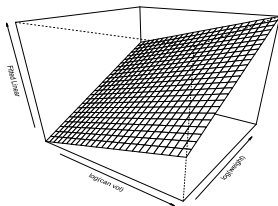
- Note: Captures interaction terms between x_1 and x_2

©Emily Fox 2014

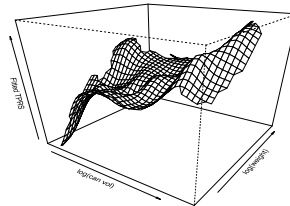
12

Tensor Product Splines Example

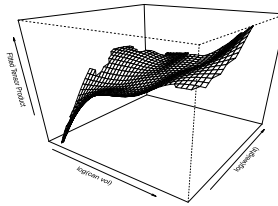
- For prostate cancer dataset, fits of log PSA as a function of log cancer volume and log weight for various models



Linear fit



Thin plate regression spline



Tensor product spline

From Wakefield textbook

similar

©Emily Fox 2014

13

Generalized Additive Models

- Both for computational reasons and added interpretability, models that assume an additive structure are very popular
- Assuming a GLM framework:

$$g(\mu(x)) = \alpha + f_1(x_1) + \dots + f_d(x_d)$$

$$LM: y = \alpha + f_1(x_1) + f_2(x_2) + \dots + f_d(x_d)$$

- Is this model identifiable?

No, can shift α and shift to compensate \rightarrow exactly same $g(\mu)$

- Can model $f_j(x_j)$ using any smoother

Fix: Constrain $\sum_{j=1}^d f_j(x_{(j)}) = 0$ to match change
many choices! (spline, kernel method, etc.) $\rightarrow \hat{\alpha} = \bar{y}$
(module 2)

©Emily Fox 2014

14

GAM Example

$$GLM: g(\mu) = X^T \beta$$

- Consider using a penalized regression spline of order p_j with L_j knots for each covariate x_j

or $g(\mu) = \beta_0 + \sum_{j=1}^d \left[\sum_{k=1}^{p_j} \beta_{jk} x_j^k + \sum_{\ell=1}^{L_j} b_{j\ell} (x_j - \xi_{j\ell})_+^{p_j} \right] = f_j(x_j)$

- Penalization is applied to the spline coefficients b_j

$$\sum_{j=1}^d \lambda_j \sum_{\ell=1}^{L_j} b_{j\ell}^2$$

Comments:

- The GAM is very interpretable
 - $f_j(x_j)$ is not influenced by the other $f_i(x_i)$
 - Can plot f_j to straightforwardly see the relationship between x_j and y
- Will see that this also leads to computational efficiencies

©Emily Fox 2014

15

Backfitting

- To begin, assume a standard (non-GLM) regression setting

$$y = f(x) + \varepsilon$$

- For concreteness, consider

$$\min_{f_1, \dots, f_d} \sum_{i=1}^n (y_i - \alpha - \sum_{j=1}^d f_j(x_{ij}))^2 + \sum_{j=1}^d \lambda_j \int f_j'(t)^2 dt$$

- Result is an **additive cubic spline model** with knots at the unique values of x_{ij}

- For X full column rank, can show that solution is unique. Otherwise, linear part of $f_j(x_j)$ is not uniquely determined

- Here, clearly $\hat{\alpha} = \bar{y} \quad \left(\sum_i f_j(x_{ij}) = 0 \right)$

- How do we think about fitting the other parameters??

©Emily Fox 2014

16

Backfitting

- **Backfitting** is an iterative fitting procedure
- Since $f(x)$ is additive, if we condition on the fit of all other components $f_j(x_j), j \neq i$, then we know how to fit $f_i(x_i)$
- Iterate the estimation procedure until convergence

Backfitting Algorithm

Algorithm 9.1 *The Backfitting Algorithm for Additive Models.*

1. Initialize: $\hat{\alpha} = \frac{1}{N} \sum_{i=1}^N y_i, \hat{f}_j \equiv 0, \forall i, j$.
2. Cycle: $j = 1, 2, \dots, p, \dots, 1, 2, \dots, p, \dots$,

$$\hat{f}_j \leftarrow \mathcal{S}_j \left[\{y_i - \hat{\alpha} - \sum_{k \neq j} \hat{f}_k(x_{ik})\}_1^N \right],$$

$$\hat{f}_j \leftarrow \hat{f}_j - \frac{1}{N} \sum_{i=1}^N \hat{f}_j(x_{ij}).$$

until the functions \hat{f}_j change less than a prespecified threshold.

From Hastie, Tibshirani, Friedman book

GAMs and Logistic Regression

- A generalized additive logistic regression model has the form
- The functions f_1, \dots, f_d can be estimated using a backfitting algorithm, too
- First, recall IRLS algorithm for *parametric* logistic regression

$$z = X\beta^{\text{old}} + W^{-1}(y - p)$$

$$\beta^{\text{new}} \leftarrow \arg \min_{\beta} (z - X\beta)^T W (z - X\beta)$$

©Emily Fox 2014

19

GAMs and Logistic Regression

Algorithm 9.2 *Local Scoring Algorithm for the Additive Logistic Regression Model.*

1. Compute starting values: $\hat{\alpha} = \log[\bar{y}/(1 - \bar{y})]$, where $\bar{y} = \text{ave}(y_i)$, the sample proportion of ones, and set $\hat{f}_j \equiv 0 \forall j$.

2. Define $\hat{\eta}_i = \hat{\alpha} + \sum_j \hat{f}_j(x_{ij})$ and $\hat{p}_i = 1/[1 + \exp(-\hat{\eta}_i)]$.

Iterate:

- (a) Construct the working target variable

$$z_i = \hat{\eta}_i + \frac{(y_i - \hat{p}_i)}{\hat{p}_i(1 - \hat{p}_i)}$$

- (b) Construct weights $w_i = \hat{p}_i(1 - \hat{p}_i)$

- (c) Fit an additive model to the targets z_i with weights w_i , using a weighted backfitting algorithm. This gives new estimates $\hat{\alpha}, \hat{f}_j, \forall j$

3. Continue step 2. until the change in the functions falls below a pre-specified threshold.

From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

20

GAM Logistic Example

- Example: *predicting spam*
- Data from UCI repository
- Response variable: *email* or *spam*
- 57 predictors:
 - 48 quantitative – percentage of words in email that match a give word such as “business”, “address”, “internet”,...
 - 6 quantitative – percentage of characters in the email that match a given character (; , [! \$ #)
 - The average length of uninterrupted capital letters: CAPAVE
 - The length of the longest uninterrupted sequence of capital letters: CAPMAX
 - The sum of the length of uninterrupted sequences of capital letters: CAPTOT

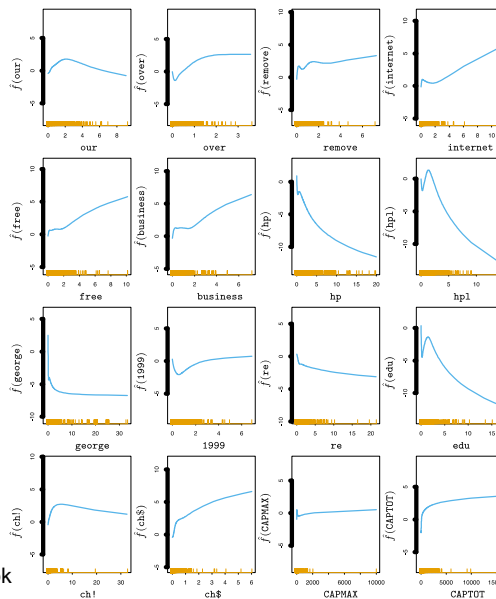
©Emily Fox 2014

21

GAM Logistic Example

- Test set of 1536 emails
- Training set: n=3065
- Use a GAM with a cubic smoothing spline
 - Each with 4 dof
- Estimated functions for significant predictors
 - Note large discontinuity near 0 for many
- Test error of 6.6%

From Hastie, Tibshirani, Friedman book



©Emily Fox 2014

22

Other GAM formulations

- Semiparametric models:

$$g(\mu) =$$

- ANOVA decompositions:

$$f(x) =$$

Choice of:

- ☐ Maximum order of interaction
- ☐ Which terms to include
- ☐ What representation

- Tradeoff between full model and decomposed model

©Emily Fox 2014

23

Connection with Thin Plate Splines

- Recall formulation that lead to natural thin plate splines:

$$\min_f \sum_{i=1}^n \{y_i - f(x_i)\}^2 + \lambda J(f)$$
$$J(f) = \int \int_{\mathbb{R}^2} \left[\left(\frac{\partial^2 f(x)}{\partial x_1^2} \right)^2 + 2 \left(\frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 + \left(\frac{\partial^2 f(x)}{\partial x_2^2} \right)^2 \right] dx_1 dx_2$$

- There exists a $J(f)$ such that the solution has the form
- However, it is more natural to just assume this form and apply

$$J(f) = J(f_1 + f_2 + \cdots + f_d) = \sum_{j=1}^d \int f_j''(t_j)^2 dt_j$$

©Emily Fox 2014

24

What you need to know

- Nothing is conceptually hard about multivariate x
- In practice, nonparametric methods struggle from curse of dimensionality
- Options considered:
 - Thin plate splines
 - Tensor product splines
 - Generalized additive models
 - Combinations (to model some interaction terms)

©Emily Fox 2014

25

Readings

- Wakefield – 12.1-12.3
- Hastie, Tibshirani, Friedman – 5.7, 9.1
- Wasserman – 4.5, 5.12

©Emily Fox 2014

26