

Module 5 Part 2: Penalized Splines

Wednesday, November 1, 2023 14:21



BIOS526_M
5_PartII_...

Module 5, part II: Penalized and Smoothing Splines

BIOS 526

1/62

Module 5, part II: Penalized and Smoothing Splines

Reading

- Sections 5.4 and 5.5 in Hastie et al.
- Sections 3.1 - 3.14, 4.9 in Ruppert et al.

Concepts

- Constraints and penalized regression.
- Smoothing matrix and smoothing parameter.
- Generalized cross-validation to choose roughness penalty.
- Mixed models to choose roughness penalty.

2/62

Module 5, part II: Penalized and Smoothing Splines

Motivating Example: Daily Temperature and Deaths

- alldeaths: daily non-accidental deaths in the 5-county New York City, 2001-2005.
- Temp: daily temperature in Fahrenheit.

```
> load ("NYC.RData")
> plot(alldeaths$Temp,xlab="Temperature (F)",ylab ="Death count",data=health)
```

3/62

Module 5, part II: Penalized and Smoothing Splines

Regression Problem

Let y_i be the number of non-accidental deaths on day i and x_i be the same-day temperature.

We consider the nonparametric regression problem:

$$y_i = g(x_i) + \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2).$$

We can approximate $g(x_i)$ using a *generalized additive model*

$$y_i = g(x_i) = \beta_0 + \beta_1 x_i + \sum_{m=1}^M \beta_{m+1} b_m(x_i) \quad \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$$

Specify $b_m(x_i)$. E.g., linear spline with 9 equidistant interior knots $\kappa_1, \kappa_2, \dots, \kappa_9$ within the observed range of daily temperature, a piecewise linear spline model is

$$g(x_i) = \beta_0 + \beta_1 x_i + \beta_2 (x_i - \kappa_1)_+ + \beta_3 (x_i - \kappa_2)_+ + \dots + \beta_{10} (x_i - \kappa_9)_+$$

4/62

Module 5, part II: Penalized and Smoothing Splines

Automatic Knot Selection

What if we don't know the number and locations of the knots?

Approach:

- Start with **a lot of knots**. This ensures that we will not miss important fine-scale behaviour.
- Assume most of the knots are not useful and **shrink** their coefficients toward zero.
- Determine how much to shrink based on some criteria (e.g. GCV or AIC).

Benefits:

- Knot placement is not important if the number is dense enough.
- Shrinking most coefficients to zero will stabilize model estimation similar to performing variable selection.

Module 5, part B: Fitting and Smoothing Splines

will limit β^* so not too big → prevent overfitting & collinearity

Penalized Spline

Consider the basis expansion:

$$g(x_i) = \beta_0 + \beta_1 x_i + \sum_{m=1}^M \beta_{m+1} b_m(x_i). \quad (1)$$

Constrain the magnitude of the coefficients β_j :

Consider the **ridge-regression** penalty:

$$\beta_0^2 + \beta_1^2 + \dots + \beta_{M+1}^2 \leq C, \quad (2)$$

equivalently,

$$\|\beta^*\|_2^2 \leq C,$$

where $\beta^* = [\beta_2, \dots, \beta_{M+1}]^T$ and C is a positive constant.

Module 5, part B: Fitting and Smoothing Splines

Lasso & Ridge

$\|L\|$ Penalties

- Ridge regression: $\|L\|$ penalty = $\|\beta^*\|_2^2$
- Other penalties: lasso = absolute value = $\|L\|$ -penalty = $\|\beta^*\|_1 = \sum_{j=1}^{M+1} |\beta_j|$
- Ridge shrinks coefficients of vectors in b-spline basis, but does not induce sparsity.
- Ridge is easy to solve - closed form solution!
- Lasso tends to make some coefficients exactly zero. Trickier to solve. More on this later in the course.
- A small C will shrink more coefficients, as well as shrink them closer to zero.
- Our goal: convert the two problems of **how many knots** and **where** to put them into a **single parameter** that we can choose. w/ GCV

Module 5, part B: Fitting and Smoothing Splines

Matrix of $g()$

For simplicity, consider a linear spline. Then evaluate the basis functions at each x_i , $i = 1, \dots, n$:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & (x_1 - \kappa_1)_+ & \dots & (x_1 - \kappa_{M+1})_+ \\ 1 & x_2 & (x_2 - \kappa_1)_+ & \dots & (x_2 - \kappa_{M+1})_+ \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & (x_n - \kappa_1)_+ & \dots & (x_n - \kappa_{M+1})_+ \end{bmatrix}$$

Then write $g(x_i)$, $i = 1, \dots, n$ in matrix form:

$$\mathbf{G} = \mathbf{X}\beta, \quad g(x_i) = \beta_0 + \beta_1 x_i + \sum_{m=1}^M \beta_{m+1} (x_i - \kappa_m)_+$$

Then the residuals are $\mathbf{Y} - \mathbf{G} = \mathbf{Y} - \mathbf{X}\beta$.



Ridge regression: introduces bias to decrease variance

↓
minimizes residuals AND $\lambda \times \text{slope}^2$?

Lasso regression: minimizes residuals AND $\lambda \times |\text{slope}|$

It's better to standardize variables before applying ridge, so it's squared so scale matters

→ ↑ λ , β goes all the way to 0, unlike ridge regression, where they will approach 0

Lasso can get rid of useless variables like astrological sign



Constrained formulation

If you don't want to penalize all the β 's

We define the objective function:

$$\underset{\beta}{\operatorname{argmin}} \quad (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) \text{ subject to } \beta' \mathbf{B}\beta \leq C.$$

Here, \mathbf{B} is a diagonal matrix with 0 and 1 entries selecting which coefficients are penalized:

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 0_{1 \times 40} \\ 0 & 0 & 0_{1 \times 40} \\ 0_{10 \times 1} & 0_{40 \times 1} & I_{40 \times 40} \end{bmatrix} \quad \text{then} \quad \beta' \mathbf{B}\beta = \sum_{m=2}^{M+1} \beta_m^2$$

not penalize intercept
not penalize global term
aka
 $0 \ 0 \ 0$
 $0 \ 0 \ 0$
 $0 \ 0 \ 1$

Penalized formulation

This problem can be equivalently formulated as

$$\underset{\beta}{\operatorname{argmin}} \quad (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \beta' \mathbf{B}\beta \quad (3)$$

There is a one-to-one mapping between λ and the constraint C . λ is often called the **smoothing parameter**. ↗ can be ↗ when estimating a spline ↗ 0 to ∞
↳ tuning parameter ↳ hyperparameter

note:
 $\frac{1}{2} (y - x\beta)^2$
 $d\beta$
 $\sim 2x(y - x\beta)$

Closed-form solution

Won't need to derive on exam

$$\begin{aligned} \underset{\beta}{\operatorname{argmin}} \quad & (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \beta' \mathbf{B}\beta. \\ & \sim 2\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) \\ & \mathbf{X}'\mathbf{X} + \lambda \mathbf{B} \beta = \mathbf{X}'\mathbf{Y} \\ & \beta = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}'\mathbf{Y}. \end{aligned}$$



Differentiate wrt β and set to zero:

$$\begin{aligned} -2\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) + 2\lambda \mathbf{B}\beta &= 0 \\ -\mathbf{X}'\mathbf{Y} + \mathbf{X}'\mathbf{X}\beta + \lambda \mathbf{B}\beta &= 0 \\ (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})\beta &= \mathbf{X}'\mathbf{Y} \\ \hat{\beta} &= (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}'\mathbf{Y}. \end{aligned}$$

Closed-form solution

penalized least squares, not ordinary least squares

The least squares solution is

$$\hat{\beta} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}'\mathbf{Y} \quad (4)$$

for some positive number λ . Note:

- When $\lambda = 0$, $\hat{\beta}$ becomes the ordinary least squares estimate. So no penalization is present ($C = \infty$).
- When $\lambda \rightarrow \infty$, $(\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1}$ becomes small, so $\hat{\beta}_j \rightarrow 0$ if $B_{jj} = 1$ ↗ so slope gets smaller

i.e., $\hat{\beta}_0$ and $\hat{\beta}_1$ are not penalized

but the non-linear parts of the splines (coeffs corresponding to knots)

go to 0 as $\lambda \rightarrow \infty$

Mortality and Temperature Example

Consider the death and mortality analysis. Assume 40 equidistant knots and linear splines:

$$y_i = \beta_0 + \beta_1 x_i + \sum_{m=1}^{40} \beta_{1+m} (x_i - K_m)_+$$

We will penalize $\beta_{1+m}, \dots, \beta_{M+1}$ using the \mathbf{B} matrix:

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 0_{1 \times 40} \\ 0 & 0 & 0_{1 \times 40} \\ 0_{10 \times 1} & 0_{40 \times 1} & I_{40 \times 40} \end{bmatrix}$$

$$\text{so } \beta' \mathbf{B}\beta = \sum_{m=1}^{40} \beta_{1+m}^2$$

How do we pick λ ?

Use cross-validation,
typically 10-fold cross validation.
look at the one with least residual variance aka MSE

Mortality and Temperature Example

Use cross-validation,
typically 10-fold cross validation
look at the one with least
residual variance aka MSE

Consider the death and mortality analysis. Assume 40 equidistant knots and linear splines:

$$y_i = \beta_0 + \beta_1 x_i + \sum_{m=1}^{40} \beta_{1+m} (x_i - \kappa_m)_+$$

We will penalize $\beta_{1+m}, \dots, \beta_{M+1}$ using the \mathbf{B} matrix:

$$\mathbf{B} = \begin{bmatrix} 0 & 0 & 0_{1 \times 40} \\ 0 & 0 & 0_{1 \times 40} \\ \vdots & \vdots & \vdots \\ 0_{40 \times 1} & 0_{40 \times 1} & I_{40 \times 40} \end{bmatrix}$$

so $\mathbf{P}' \mathbf{B} \mathbf{P} = \sum_{m=1}^{40} \beta_{1+m}^2$

13/52

Module 5, part B: Prediction and Smoothing Spline

Creating piecewise linear spline

We can create a design matrix with piecewise linear splines.

```
> knots = seq(range(healthTemp)[1], range(healthTemp)[2], length.out = 40+2)
> # place knots evenly on interior of the range of x
> knots = knots[c(2:(length(knots)-1))]
> X = chnderp(1, length(healthTemp), health$Temp)
> for (i in 2:(length(knots))) {
+   X[i] = colint(X, (health$Temp-knots[i])*(health$Temp>knots[i]))
+ }
> B = diag(42)
> B[1,1]=0
> B[2,2]=0
> dim (X); dim (B)
[1] 1826 42
[1] 42 42
```

14/52

Module 5, part B: Prediction and Smoothing Spline

Mortality and Temperature Example

We now search through different values of λ . For each λ , we will

- Calculate the penalized $\hat{\beta}$.
- Calculate $\hat{\beta}' \mathbf{B} \hat{\beta}$.
- Calculate the fitted value $\hat{Y} = \mathbf{X} \hat{\beta}$.
- Calculate the GCV using the matrix: $\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1}\mathbf{X}'$.

We will select the λ with the smallest GCV.

use GCV

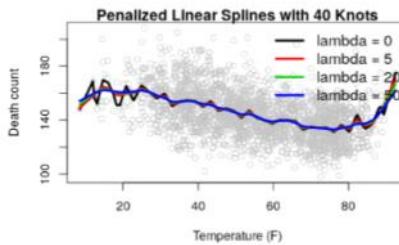
```
> Y = health$alldeaths
> lambda = 0
> beta = solve(t(X)X*X + lambda*B) X*t(X) X*Y
> H = X Z*X solve (t(X)X*X + lambda*B) X*X t(X) ##Hat matrix
> That = X*X*beta ##Fitted values
> GCV = mean ((Y-That)^2) / (1- mean (diag(H)))^2
> C = t(beta)*X'*H*X*beta
```

loop over all values
of λ in the
order, you
can plot
on plot
what
that
means
like
this

15/52

Module 5, part B: Prediction and Smoothing Spline

Effects of Penalization

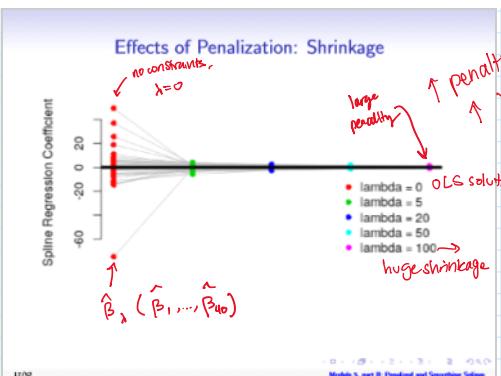


16/52

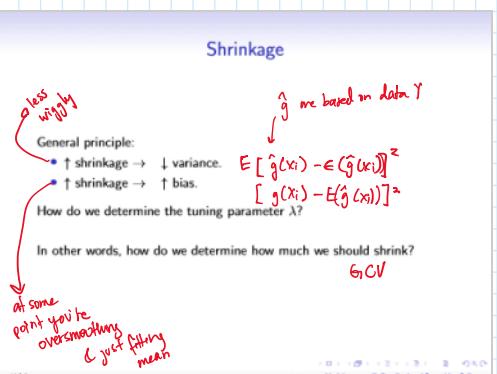
Module 5, part B: Prediction and Smoothing Spline

Higher penalty

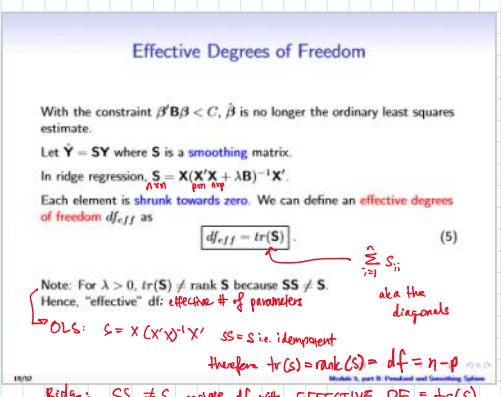
As $\lambda \uparrow$,
line becomes
less wiggly;
betas go towards 0?



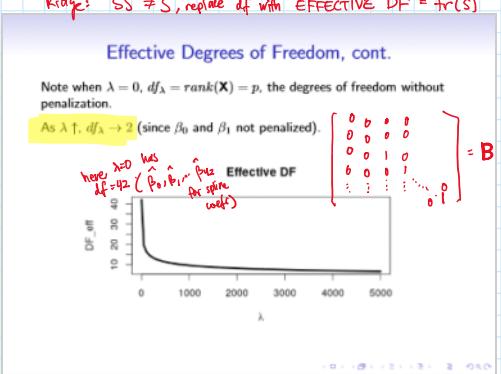
Module 5, part B: Fitted and Smoothing Spline



Module 5, part B: Fitted and Smoothing Spline



Module 5, part B: Fitted and Smoothing Spline



Module 5, part B: Fitted and Smoothing Spline

Generalized Cross-validation Error, revisited

We previously defined GCV:

$$GCV = \frac{1}{n} \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{[1 - n^{-1} \text{tr}(\frac{1}{\lambda} \mathbf{H})]^2}$$

Note that $\hat{\mathbf{Y}} = \mathbf{HY}$ where $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$.

Now we can apply GCV to **any** prediction of \mathbf{Y} that can be written in the form:

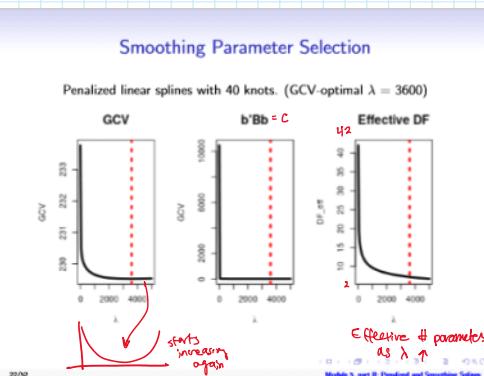
$$\hat{\mathbf{Y}} = \mathbf{SY}.$$

Then GCV is defined:

$$GCV = \frac{1}{n} \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{[1 - n^{-1} \text{tr}(\frac{1}{\lambda} \mathbf{S})]^2}$$

This is the definition we will use hereafter.

Module 5, part B: Fitted and Smoothing Spline



Residual Error Variance Estimate

Recall our model is

$$y_i = g(x_i) + \epsilon_i \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2).$$

We now have an estimate $\hat{g}(x_i)$. How about σ^2 ?

We have two options:

$$\text{mgcv:: gam}$$

$$\text{c. if } \hat{\sigma}^2 = \frac{\sum_{i=1}^n [y_i - \hat{g}(x_i)]^2}{n - df_{\text{eff}}} \quad (6)$$

The above is a biased estimate. Some software gives you the option to use

$$\hat{\sigma}^2_{\text{unbiased}} = \frac{\sum_{i=1}^n [y_i - \hat{g}(x_i)]^2}{n - 2\text{tr}(\mathbf{S}) + \text{tr}(\mathbf{SS}')}. \quad (7)$$

Ref. of
will use this
later in
mgcv package
is R

G in OLS, n-2p+p=n-p✓

Module 5, part B: Fitted and Smoothing Spline

Variance of $\hat{g}(x_i)$

Now we can calculate uncertainty associated with $\hat{g}(x_i)$ at each x_i .

With slight abuse of notation, let \mathbf{x}_i' be the row vector of basis function values for x_i .

The variance of $\hat{g}(x_i)$ is

$$\begin{aligned} \text{Var}[\hat{g}(x_i)] &= \text{Var}[\mathbf{x}_i' \hat{\beta}] = \mathbf{x}_i' \text{Var}[\hat{\beta}] \mathbf{x}_i \\ &= \mathbf{x}_i' \text{Var}[(\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}' \mathbf{Y}] \mathbf{x}_i = \tilde{\mathbf{x}}_i' [\mathbf{X}'\mathbf{X} + \lambda \mathbf{B}]^{-1} \mathbf{X}' \text{Var}(\mathbf{Y}) [\mathbf{X}'\mathbf{X} + \lambda \mathbf{B}]^{-1} \tilde{\mathbf{x}}_i \\ &= \sigma^2 \mathbf{x}_i' (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} (\mathbf{X}'\mathbf{X}) (\mathbf{X}'\mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{x}_i. \end{aligned}$$

Note: you should decide whether or not to include the variance due to the intercept. If $x_i[1] = 1$, then the variance estimate of $\hat{g}(x_i)$ includes this source of uncertainty.

Is $\hat{g}(x_i)$ including β_0 ? MGCV:: gam does not include β_0

Module 5, part B: Fitted and Smoothing Spline

Confidence interval and prediction interval

Obtain point-wise confidence interval derived from previous expression by plugging in $\hat{\sigma}^2$ for σ^2 .

If $\lambda = 0$: the previous equation reduces to the OLS variance.

Similarly the variance for an unobserved point y_i^* with covariate x_i^* has variance **new value of (y_i^*, \hat{y}_i^*)**

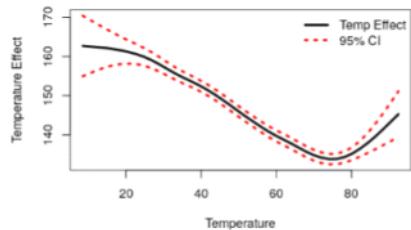
$$\text{Var}[y_i^*] = \sigma^2 + \sigma^2 x_i^* (\mathbf{X}'\mathbf{X} + \lambda\mathbf{B})^{-1} (\mathbf{X}'\mathbf{X}) (\mathbf{X}'\mathbf{X} + \lambda\mathbf{B})^{-1} x_i^*.$$

↳ same as previous slide, just add additional σ^2

26/52 Module 5, part B: Prediction and Smoothing Systems

Temperature Effect on Mortality: pointwise CI

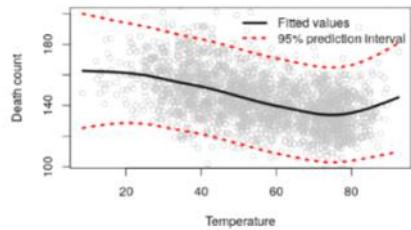
```
> Upper95.ci = That + 1.96* sqrt(diag(pred.vcov))
> Lower95.ci = That - 1.96* sqrt(diag(pred.vcov))
```



26/52 Module 5, part B: Prediction and Smoothing Systems

Daily Mortality Prediction

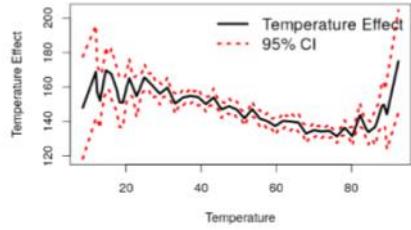
```
> Upper95 = That + 1.96* (signal + sqrt(diag(pred.vcov)))
> Lower95 = That - 1.96* (signal + sqrt(diag(pred.vcov)))
```



27/52 Module 5, part B: Prediction and Smoothing Systems

Temperature Effect on Mortality

Compare to a model without penalization ($\lambda = 0$).



28/52 Module 5, part B: Prediction and Smoothing Systems

Smoothing Splines: other penalties

NOT PART OF EXAM

A function with large second derivatives can be interpreted as rougher, as the function is allowed to change very rapidly.

We now add a "roughness" penalty to encourage smoothness:

$$\hat{g}(x) = \arg \min_{g \in \mathcal{G}} \{\mathbf{Y} - g(x)\}' \{\mathbf{Y} - g(x)\} + \lambda \int_a^b \{g''(x; \beta)\}^2 dx. \quad (8)$$

where \mathcal{G} are twice-differentiable functions, $x \in \mathbb{R}^n$ is the vector of x_i , $i = 1, \dots, n$, and a and b is the range of x .

Smoothing splines:
simplifying of
other than
 χ^2

29/52

Module 5, part B: Pseudol and Smoothing Spline

Smoothing spline, cont.

$$\hat{g}(x) = \arg \min_{g \in \mathcal{G}} \{\mathbf{Y} - g(\bar{x})\}' \{\mathbf{Y} - g(\bar{x})\} + \lambda \int \{g''(x; \beta)\}^2 dx.$$

where \mathcal{G} is the class of twice-differentiable functions and $\bar{x} \in \mathbb{R}^n$ is the vector of x_i , $i = 1, \dots, n$.

- Note that first derivatives are not penalized.
- The second part uses the squared second derivative that is a good measure of roughness.
- Shrinks coefficients in a cubic polynomial, causing function to change less quickly.
- λ determines the relative importance of minimizing the residual sum of squares or the roughness.

30/52

Module 5, part B: Pseudol and Smoothing Spline

Smoothing Spline

It turns out the solution $\hat{g}(x)$ is a "natural cubic spline" (a cubic spline with linearity at the boundaries) with knots at the observed points x_i .

More generally, the objective function in (8) with penalized second derivatives is equivalent to basis expansion of x

$$(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) + \lambda \beta' \mathbf{B} \beta \quad (9)$$

for a certain \mathbf{B} matrix based on second moments of the basis functions, no longer diagonal; see Ruppert et al p. 75. \rightarrow semiparametric regression

The key point is that (9) is a general formula applying to different ridge-like penalties for certain \mathbf{B} .

As before,

- for a given λ , we can estimate $g(x)$ using penalized least squares;
- search through λ to minimize GCV or another criterion.

31/52

Module 5, part B: Pseudol and Smoothing Spline

Note: another package called mgcv but not as good

Package mgcv in R \rightarrow WILL BE ON EXAM

The mgcv (Mixed GAM Computation Vehicle) package in R contains the `gam()` function to fit a large variety of smoothing splines with automatic smoothing parameter selection. We will examine different options throughout the class.

Default option is given in parenthesis.

- Basis functions (default: thin plate regression spline).
- Basis dimension (default: $k = 10$ with one constraint: $\sum \hat{g}(x_i) = 0$, makes max edf=9, min edf=1 (unpenalized linear term, if edf=1, results nearly equivalent to lm()). 10 parameters - 1 constraint = 9 free parameters
- Selection methods (default: GCV).
- Family (default: Gaussian).
- Standard error computation (default: Bayesian).

32/52

Module 5, part B: Pseudol and Smoothing Spline

Temperature Effect on Mortality

```
> library(mgcv)
> fit1 = gam(maldeaths~n(Temp), data= health)
> summary(fit1)

Family: gaussian
Link function: identity

Parametric coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 143.917    0.364   407 <2e-16 ***
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms: → approximate F statistic
edf Ref.df F p-value
n(Temp) 6.03 7.2 80.6 <2e-16 ***
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.241 Deviance explained = 24.3%
GCV = 229.47 Scale est. = 228.58 n = 1826
```

33/92 Module 5, part II: Pseudol and Smoothing Spline

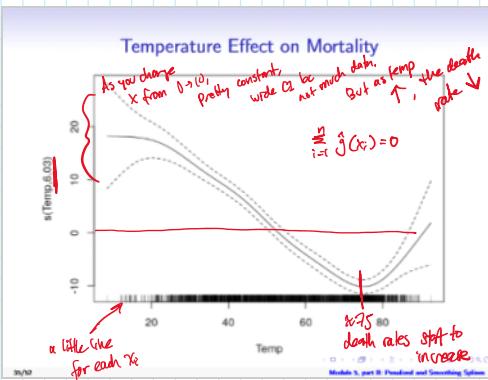
mgcv::gam output

```
Approximate significance of smooth terms:
edf Ref.df F p-value
n(Temp) 6.03 7.2 80.6 <2e-16 ***
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.241 Deviance explained = 24.3%
GCV = 229.47 Scale est. = 228.58 n = 1826
```

G ↗
• edf = effective Df for $tr(S)$.
• Ref edf = effective Df for $2tr(S) - tr(S'S)$.
• Scale est. = estimated residual error σ^2 (using edf).
• F statistic: approximate significance of Temp. Uses Ref edf.
• Use plots to interpret $\hat{j}(x_i)$.

34/92 Module 5, part II: Pseudol and Smoothing Spline



Checking gam

The default is $k = 10$, such that highest possible EDF is 9 (because of identifiability constraint).

```
> gam.check(fit1) → function we use to get diagnostic plots
```

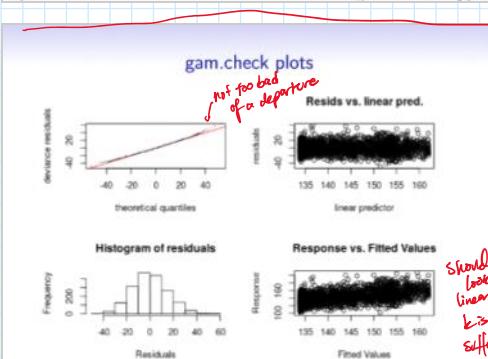
Method: GCV Optimizer: mgcv
Smoothing parameter selection converged after 5 iterations.
The RMS GCV score gradient at convergence was 7.242e-05.
The Hessian was positive definite.
Model rank = 10 / 10

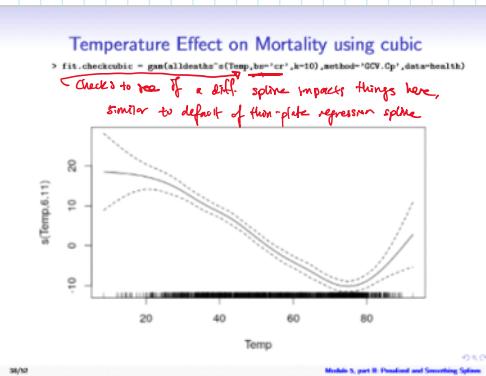
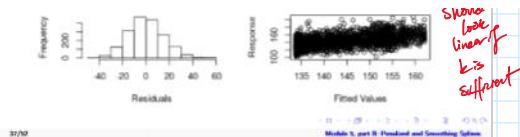
Basis dimension (k) checking results. Low p-value (k-index<1) may indicate that k is too low, especially if edf is close to k^{*}.

k^* edf k-index p-value → $H_0: \text{Basis dimension is adequate}$
 $n(Temp) 9.00 6.03 1.02 0.68$ → This test not very useful

In general, "k" should be "notably" higher than edf
e.g. edf=5, here we will check fitting model w/a higher basis dimension

36/92 Module 5, part II: Pseudol and Smoothing Spline





(Inferential sensitivity analysis to basis dimension.)

Temperature Effect on Mortality

Thin plate splines with $k = 40$.

```
> fit2=gam(alideaths~s(Temp, k = 40), data = health)
> summary(fit2)
```

Formula:
 $\text{alideaths} \sim s(\text{Temp}, k = 40)$

Parametric coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	143.917	0.354	407	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms:

edf	Ref.df	F p-value
s(Temp)	6.23	7.86 73.9 <2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.241 Deviance explained = 24.3%
 GCV = 229.51 Scale est. = 228.6 n = 1826

Extract Useful Model Statistics

Full list see ?gamObject.

- AIC (with edf at penalized estimates)


```
> AIC (fit)
[1] 15109.62
```
- Variance-covariance matrix of the coefficients of your spline basis


```
> dim (fit$Ve) ### Frequentist's
[1] 10 10
> dim (fit$Vp) ### Bayesian
[1] 10 10
```

(not really; they have frequentist coverage probabilities, so we default to using these)
- Fitted value


```
> fit$fitted
```

A different approach to smoothing rather than GCV

Penalized splines as BLUPs

- GCV may undersmooth.
- An alternative is to treat the coefficients of the truncated polynomials as random effects, and then use BLUPs.
- For concreteness, consider a linear spline:

Applies to non-differentiable data (though you could use it with?)

$y_i = \beta_0 + \beta_1 x_i + \sum_{m=1}^M \theta_m (x_i - \kappa_m)_+ + \epsilon_i$

$\theta_m \stackrel{\text{iid}}{\sim} N(0, \tau^2)$, $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$

random slope

Data is not longitudinal here, so this is little trick

\uparrow wiggly $\uparrow \tau^2$
 \uparrow noisy $\uparrow \sigma^2$

θ_m becomes larger with more wiggly data
 ϵ_i becomes larger with more noisy data

$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_M \end{bmatrix}$ $\Theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_M \end{bmatrix}$ $\mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$ $\mathbf{Z} = \begin{bmatrix} (x_1 - \kappa_1)_+ & \dots & (x_1 - \kappa_M)_+ \\ \vdots & \ddots & \vdots \\ (x_n - \kappa_1)_+ & \dots & (x_n - \kappa_M)_+ \end{bmatrix}$

Mixed model for estimating a penalized spline

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_1 & \dots & X_n \\ K_1 & \dots & K_n \end{bmatrix} \beta + \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Given τ^2 and σ^2 , we seek to minimize

$$\frac{1}{\sigma^2} \|Y - X\beta - Z\Theta\|^2 + \frac{1}{\tau^2} \|\Theta\|_2^2$$

which we can think of ridge regression with penalty $\lambda = \frac{\tau^2}{\sigma^2}$.

We estimate all parameters from the data using the mixed modeling tools we previously learned, and thus obtain a model-based estimate of λ .

$$\|Y - X\beta - Z\Theta\|^2 + \frac{\tau^2}{\sigma^2} \|\Theta\|_2^2$$

If τ^2 is very large \rightarrow large penalty
 \rightarrow shrink coefficients towards 0

40/52

Module 9, part B: Prediction and Smoothing Systems

Selecting penalty using mixed models

- In mgcv::gam, we can use the option method='REML'.
- Often results in greater smoothing (smaller effective degrees of freedom)

```
> fit.real = gam(alldaths~n(Temp,bs="tp",k=10),method="REML", data= health)
> summary(fit.real)

Family: gaussian
Link function: identity

Formula:
alldaths ~ n(Temp, bs = "tp", k = 10)

Parametric coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 143.9168 0.3639 406.7 <2e-16 ***
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Approximate significance of smooth terms:
df Ref.df F p-value
n(Temp) 5.499 6.665 86.06 <2e-16 ***
...
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

R-sq.(adj) = 0.24  Deviance explained = 24.3%
-RMSE = 7555.7 Scale est. = 238.66 n = 1826
```

41/52

Module 9, part B: Prediction and Smoothing Systems

Estimate the slope at a particular x_i

In linear regression $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$.

In GAMs, we have $\hat{y}_i = \hat{\beta}_0 + \hat{g}(x_i)$, and slope changes with x_i .

What is the rate of change at 40 degrees Fahrenheit?

```
> # I usually check whether this is consistent with the plot
> new <- health[, 1] # grab any row, we are going to change temperature only
> newTemp <- 40 - 1e-06 # subtract some small number
> y1 <- predict(fit.real,new)
> newTemp <- 40 + 1e-06 # add some small number
> y2 <- predict(fit.real,new)
> y1/y2 - 1
49
-0.526
```

42/52

Module 9, part B: Prediction and Smoothing Systems

Interpretation

We interpret smoothers $\hat{g}(x_i)$ by looking at plots.

$\hat{g}'(40)$

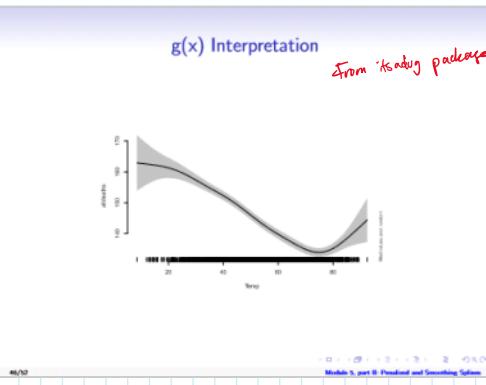
We can add some details regarding the slopes at particular x_i .

Deaths are highest at cold temperatures (< 10 degrees F) and slightly decreasing until approximately 20 degrees. Then deaths decrease at a similar rate from approximately 25 to 75 degrees. The number of deaths decreases by approximately 0.5 people / degree in a neighborhood of 40 degrees. Then the number of deaths starts to increase around 75 degrees. At 85 degrees, the number of deaths increases by approximately 0.8 for every 1 degree increase in temperature.

$\hat{g}'(85)$

43/52

Module 9, part B: Prediction and Smoothing Systems



Additive model with random intercept

GRAM: generalized additive mixed model

Recall the Nepal arm circumference dataset.

Data on 200 children collected at a maximum of 5 time points about 4 months apart.

Consider a non-linear effect of age and a random intercept:

$$\begin{aligned} arm_{ij} &= \beta_0 + g(agg_{ij}) + \theta_i + \epsilon_{ij} \\ \theta_i &\stackrel{iid}{\sim} N(0, \tau^2) \\ \epsilon_{ij} &\stackrel{iid}{\sim} N(0, \sigma^2) \end{aligned}$$

$$\theta_i \perp e_i$$

Additive model with random intercept

Syntax for (if id)

```

fit_gmm = gausmix((xgpl)@(id,1), mu = "xc"), method = "EM", data = xgpl)
par(g oma = c(0,0,0,0))
plot(xgpl, fit_gmm, main = "GMM Fit", xlab = "X", ylab = "Y")

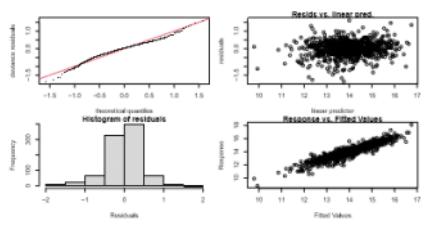
```

Basic dimensionality checking results. Low p-value (k-indep) suggests

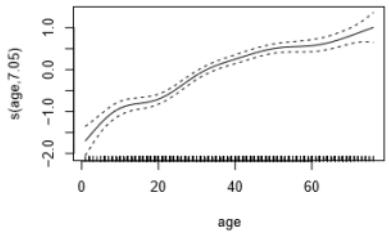
- MGCV tests random effects in an annoying way
- I tend to prefer REML

- I tend to prefer REML *but there's actually* σ^2 from R
 - EDF somewhat close to k' . Other diagnostics okay.
 - R code looks at $k = 20$ and results are similar (edf=8.5), so either this model or the one with $k = 20$ is fine.
 - For random effect, k' equals number of subjects. Good to check this because will try to fit a smooth if you don't code as a factor.

Additive mixed model with random intercept



Effect of age on arm circumference



Model S, part B: Fitted and Smoothing Spline

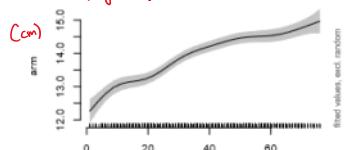
50/52

Effect of age on arm circumfencerece

This plot includes the intercept:

```
> library(mgcv)
> plot_smooth(fit.gam,view="age",rm.ranef=TRUE)
Summary:
* age : numeric predictor; with 30 values ranging from 1.000000 to 76.000000.
* id : factor; set to the value(s): 3. (Might be canceled as random effect, check below.)
* NOTE : The following random effects columns are canceled: s(id)
```

Estimate of population trajectory *no random effect?*



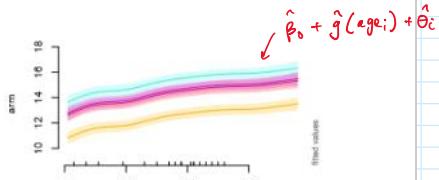
Model S, part B: Fitted and Smoothing Spline

51/52

Effect of age on arm circumference

We can also plot a few of the curves+random effects.

```
myline=c(0,18)
plot_smooth(fit.gam,real,view="age",cond=list(id=10),col="orange",ylim=myline)
plot_smooth(fit.gam,real,view="age",cond=list(id=40),col="red",add=TRUE,ylim=myline)
plot_smooth(fit.gam,real,view="age",cond=list(id=120),col="purple",add=TRUE,ylim=myline)
plot_smooth(fit.gam,real,view="age",cond=list(id=90),col="turquoise",add=TRUE,ylim=myline)
```



Model S, part B: Fitted and Smoothing Spline

52/52