

Module 5 Part 3: GAMs

Wednesday, November 8, 2023 14:15



BIOS526_M
5_PartIII_...

Module 5: Generalized Additive Models

[illegible]

Generalized Additive Models

Multiple Smoothers

Elasticity Options

Generalized Additive Models

GAMs allow us to model multiple NONLINEAR predictors

Module 3: Generalized Additive Models

Generalized Additive Models

Binary Outcome Example

Recall the example from Module 4:

Dataset: a cohort of live births from Georgia born in the year 2001 ($N = 77,340$).

Variables:

- ptb*: indicator for whether the baby from pregnancy i was born preterm (< 37 weeks).
- age*: the mother's age at delivery.
- male*: indicator of the baby's sex (1 = male; 0 = female).
- tobacco*: indicator for mother's tobacco use during pregnancy (1 = yes; 0 = no)

4/50

Generalized Additive Models

Previous analysis

```
## Fit logistic regression model
> fit = glm(ptb~age + male+tobacco, data = dat, family = binomial(link='logit'))
> summary(fit)
```

Call:
glm(formula = ptb ~ age + male + tobacco, family = binomial(link = "logit"), data = dat)

Deviance Residuals:

	1Q	Median	3Q	Max
Deviance Residuals	-0.5159663	-0.4235975	-0.4102807	2.2499946

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.421266448	0.063135527	-38.36030	< 2.2e-16 ***
age	-0.000629473	0.002159576	-0.29148	0.7706943
male#	0.072365862	0.025867177	2.79759	0.0051485 **
tobacco	0.409649486	0.063662733	7.06234	1.8258e-14 ***

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 44907.564 on 77339 degrees of freedom
Residual deviance: 44845.557 on 77336 degrees of freedom
AIC: 44853.557

1/50

Generalized Additive Models

The plot diagnostics are not very helpful with binary responses:

6/50

Generalized Additive Models

Generalized Additive Model

To account for non-linear age effect

Note: Sometimes called "semi-parametric" if including smooths & linear terms together

$$ptb_i \stackrel{\text{IND}}{\sim} \text{Bernoulli}(p_i)$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{male}_i + \beta_2 \text{tobacco}_i + s(\text{age}_i),$$

No ERROR TERM

where $s(\text{age}_i)$ is a smooth function of age. The above model is known as a **generalized additive model**.

GAMs: Generalized Additive Models.

Everything we learned about additive models is directly applicable in this setting!

7/50

Generalized Additive Models

GAM with logit link

Using the `mgcv::gam` function:

GAMs: framework

A multiple linear regression model is written as:

$$y_i = \beta_0 + \underbrace{\beta_1 x_{i1}}_{\text{linear terms}} + \underbrace{\beta_2 x_{i2}}_{\text{linear terms}} + \cdots + \underbrace{\beta_p x_{ip}}_{\text{linear terms}} + \epsilon_i$$

In a generalized additive model (GAM), we replace all the linear terms with arbitrary functions f_j :

$$y_i = \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \cdots + f_p(x_{ip}) + \epsilon_i$$

e.g., $y_i = \beta_0 + \beta_1 \log(x_{i1}) + \beta_2 \text{sqrt}(x_{i2}) + \beta_3 \log(x_{i3}) + \epsilon_i$

How do we (do we even need to) guess the exact forms (e.g., log, sqrt) for the f_j 's? We actually don't need to pre-specify!

Generalized Additive Models

GAM with logit link

Using the `mgcv::gam` function:

```
> fit.gam = gam(pfb ~ s(age) + male + tobacco, family = binomial, data = dat)
> summary(fit.gam)
```

Family: binomial
Link function: logit

Formula: *NOTE: actually fitting logit of expected value of preterm births: $\text{logit}[E(\text{pfb}_i)] \sim s(\text{age}_i) + \text{male}_i + \text{tobacco}_i$*

ptb ~ s(age) + male + tobacco

Parametric coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.44226	0.01913	-127.647	< 2e-16 ***
maleM	0.07274	0.02588	2.811	0.00494 **
tobacco	0.39016	0.05356	7.284	3.24e-13 ***

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

Approximate significance of smooth terms:

	edf	Ref.df	Chi.sq	p-value
s(age)	3.314	4.146	70.17	3.89e-14 ***

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

R-sq.(adj) = 0.00177 Deviance explained = 0.304%
URRE = -0.42095 Scale est. = 1 n = 77340

Generalized Additive Models

GAM diagnostics

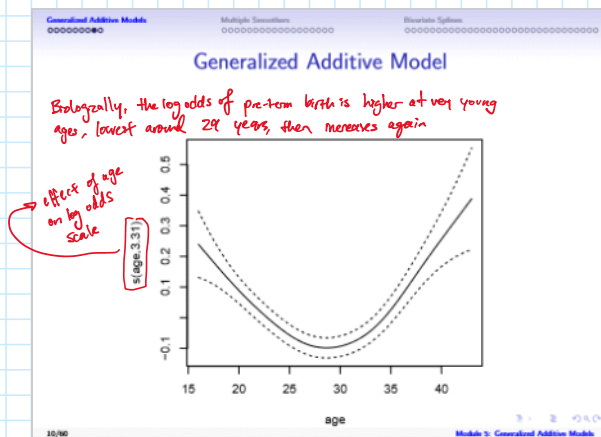
```
> gam.check(fit.gam)
```

Method: URRE *→ unbiased risk estimation replaces GCV when using a glm likelihood*
optimizer: outer newton
full convergence after 3 iterations.
Gradient range [9.816712095e-07, 9.816712095e-07]
(score -0.420948606 & scale 1).
Hessian positive definite, eigenvalue range [1.285936713e-06, 1.285936713e-06].
Model rank = 12 / 12

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
s(age)	9.00	3.31	0.94	0.66

3.31 < 9, suggests k=10 is sufficient



Generalized Additive Models

Bayesian credible intervals

We are assuming the underlying function is smooth. We can formalize this as a prior in a Bayesian model. We won't get into details: ~~return to this topic in M7~~

In Bayesian statistics, the parameters β are random variables. A ridge penalty corresponds to an improper Gaussian prior.

$$f_{\beta} \propto \exp(-\beta' B \beta / 2).$$

ridge penalty = gaussian prior

For Gaussian data, this results in a posterior distribution

$$[\beta | y, \lambda] \sim N \left\{ \hat{\beta}, \sigma^2 (X'X + \lambda B)^{-1} \right\}$$

→ in mgcv, estimated using GCV or REML-trick

For a general likelihood, we use the Fisher Information matrix (Hessian of the negative log likelihood at $\hat{\beta}$)

$$[\beta | y, \lambda] \sim N \left\{ \hat{\beta}, (\hat{I} + \lambda B)^{-1} \right\}$$

Multiply by ϕ for quasi-Poisson.

In particular, the "Bayesian credible intervals" plotted in `mgcv::gam` have frequentist coverage probabilities. See Section 6.10 in Wood.

Generalized Additive Models
○○○○○○○○○

Multiple Smoother Terms
●○○○○○○○○○○○○○○○○○○○○

Multiple Smooth Terms

Let's consider an additive model for two continuous variables, x_i and z_i :

$$y_i = \beta_0 + g_1(x_i) + g_2(z_i) + \epsilon_i, \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad (1)$$

where $g_1(\cdot)$ and $g_2(\cdot)$ denote smooth relationships between the response y_i and predictors x_i and z_i .

Again, extends to generalized linear models (binomial, Poisson, etc).

We again express non-linear functions using basis functions:

$$g_1(x_i) = \sum_{m=1}^{M_1} \beta_{m1} b_{m1}(x_i) \qquad g_2(z_i) = \sum_{m=1}^{M_2} \beta_{m2} b_{m2}(z_i)$$

Induce smoothing by penalizing regression coefficients.

Note we also use smoothers to control for **confounders** flexibly.

13/90

Generalized Addition Models
000000000000

Multiple Smoothers
000000000000000000000000

Bivariate Scales
000000000000000000000000

Associations between Mortality and Fine Particulate Matter

Fine particulate matter ($PM_{2.5}$):

- represents a mixture of solid and liquid particles in the air that are less than $2.5 \mu m$ in diameter;
- mainly arises from combustion sources (power generation, vehicle, and industrial operations).

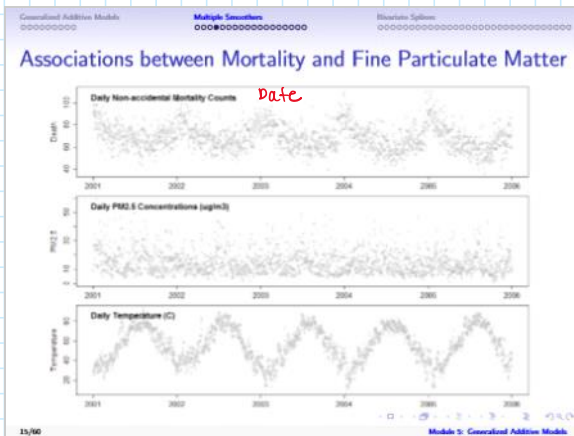
Scientific Question: what is the association between daily mortality counts and daily concentration of outdoor $PM_{2.5}$ air pollution?

Data Sources:

- Daily counts of non-accidental deaths ($age \geq 65$) in the 5 county New York City area (2001-2005) obtained from the National Center for Health Statistics (CDC).
- Daily $PM_{2.5}$ concentrations from Environmental Protection Agency
- Daily meteorology conditions from the National Climatic Data Center (NOAA).

14/90

Module 1: Generalized Addition Models



Generalized Additive Models Multiple Smoothing Bivariate Splines

Time-Series Health Model

- We are interested in the association between **daily variation** in mortality counts and **daily variation** in exposure. \rightarrow $PM_{2.5}$
- In a time-series design we view population as the unit of analysis. I.e. outcome = total mortality counts arising from the population.
- Confounders that vary smoothly in time can be easily controlled for by including smooths.

16/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Time-Series Health Model

Let y_t denote the death count on day t , and x_t be the corresponding $PM_{2.5}$ level.

There is typically a *temporal delay* between exposure and outcome.

Let's examine the association between daily mortality and **previous-day** exposure.

\rightarrow $\log y_t = \beta_0 + g_1(x_{t-1}) + \text{confounders} + \epsilon_t, \epsilon_t \stackrel{iid}{\sim} N(0, \sigma^2). \quad (2)$

We are interested in $g_1(x_{t-1})$, the smooth of lagged $PM_{2.5}$. Here we model log-transformed death counts, which can improve normality of residuals.

Confounders to consider:

- Day of the week.
- Seasonality and long-term trends.
- Same-day temperature.
- Same-day **dew point** temperature.
- Previous day's temperature and dew point temperature.

17/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Data

```
> load("NYC.RData")
> str(health)
'data.frame': 1826 obs. of 12 variables:
 $ date      : Date, format: "2001-01-01" "2001-01-02" "2001-01-03" ...
 $ alldeaths : int 171 198 179 169 201 182 167 167 193 159 ...
 $ age65plus : int 122 146 133 128 145 141 126 116 142 124 ...
 $ cardioresp: int 103 106 109 90 120 101 102 101 115 101 ...
 $ cr65plus  : int 90 92 95 77 98 90 88 82 92 88 ...
 $ dow       : chr "Monday" "Tuesday" "Wednesday" "Thursday" ...
 $ pm25      : num [1:1826, 1] 8.72 13.39 22.9 26.76 25.89 ...
 ... attr(*, "dimnames")=List of 2
 .. $ : chr "1" "2" "3" "4" ...
 .. $ : NULL
 $ Temp      : num 27.6 25.1 25.3 29.8 29.8 33.9 34.9 35.7 32.3 27.6 ...
 $ DpTemp    : num 14.2 11.8 13.1 15.9 20.5 26.7 22.2 31.2 24.4 11.7 ...
 $ rnTemp    : num 27.7 26.8 26 26.7 28.3 ...
 $ rnDpTemp  : num 17.5 14.3 13 13.6 16.5 ...
```

- $DpTemp$ = Dew point temperature.
- rm denotes 3-day **running mean** of the current day and 2 days prior.
- dow (day of week) is recorded as character.

18/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Lagged $PM_{2.5}$ Exposure

\rightarrow use default k need higher k bc effects of data are very non-linear

```
> fit = gam(log(cr65plus) ~ s(pm25.lag1) + fdow + s(date2, k = 100) + s(Temp) +
+ s(DpTemp) + s(rnTemp) + s(rnDpTemp), data = health)
> summary(fit)
```

Family: gaussian
Link function: identity

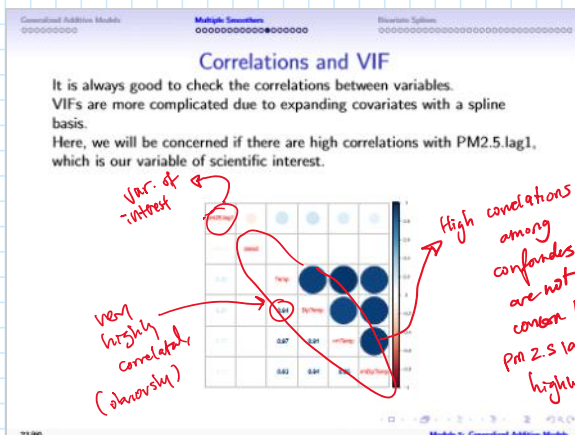
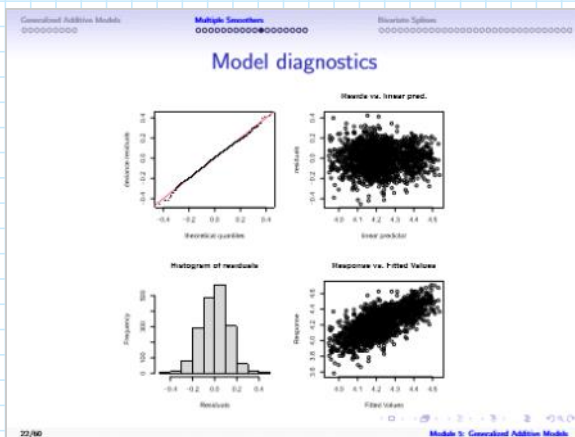
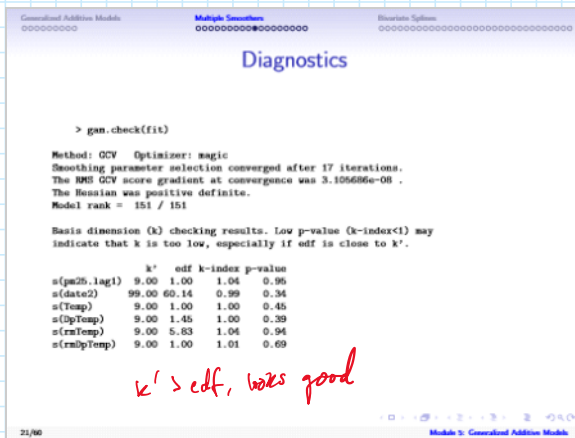
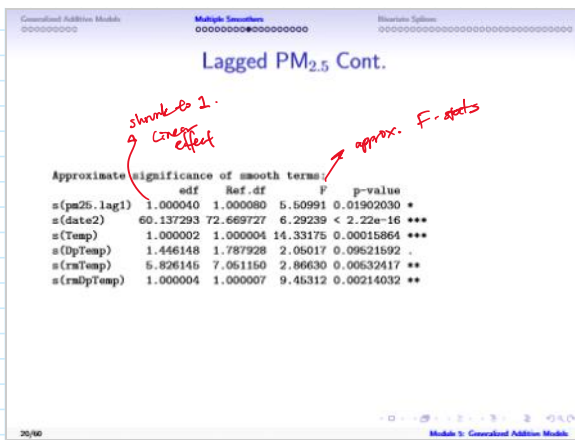
Formula:
 $\log(\text{cr65plus}) \sim s(\text{pm25.lag1}) + \text{fdow} + s(\text{date2}, k = 100) + s(\text{Temp}) + s(\text{DpTemp}) + s(\text{rnTemp}) + s(\text{rnDpTemp})$

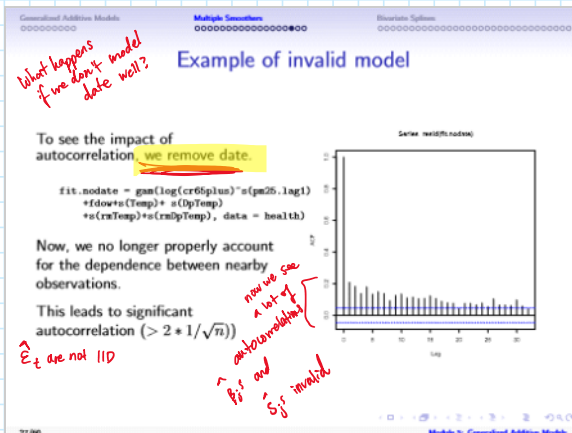
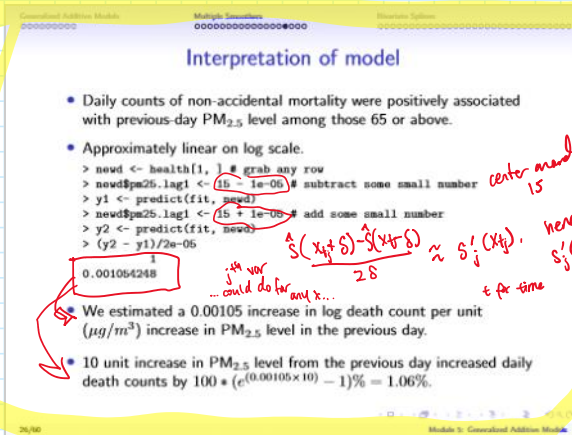
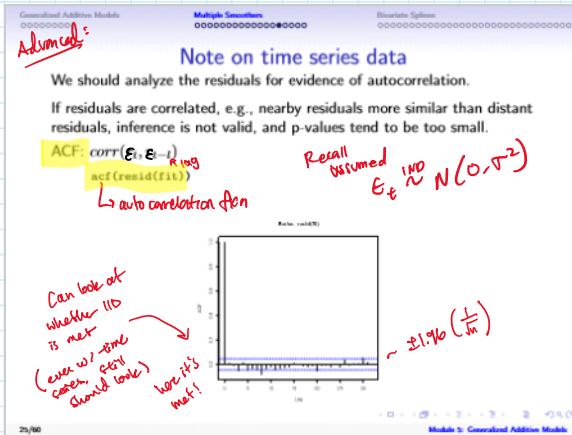
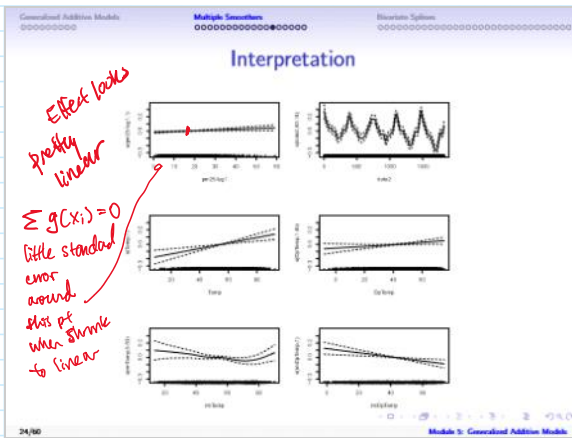
Parametric coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.265748765	0.007731497	543.97601	< 2.2e-16 ***
fdowMonday	0.036176289	0.010980266	3.29466	0.0010051 **
fdowSaturday	0.004256722	0.010918397	0.38987	0.6966825

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

19/60 Module 5: Generalized Additive Models





Generalized Additive Models Multiple Smoothers Bivariate Splines

Example of invalid model, continued

```
> summary(fit.mdate)
```

Family: gaussian
Link function: identity

Formula:
log(cr60plus) ~ s(pm25.lag1) + fdow + s(Temp) + s(DpTemp) + s(rnDpTemp)

Parametric coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.2196867	0.0086125	489.950	<2e-16 ***
fdowFriday	-0.0149662	0.0121799	-1.229	0.2193

Approximate significance of smooth terms:

	edf	Ref.df	F	p-value
s(pm25.lag1)	1.000	1.000	47.599	<2e-16 ***
s(Temp)	1.483	2.177	4.490	0.0403 *
s(DpTemp)	1.000	1.000	1.406	0.2359
s(rnDpTemp)	7.705	8.572	6.781	<2e-16 ***
s(rnDpTemp)	2.111	2.762	1.906	0.1233

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

R-sq. (adj) = 0.337 Deviance explained = 34.4%
GCV = 0.019388 Scale est. = 0.01917 n = 1825

26/60 Module 5: Generalized Additive Models

*invalid bc inflated Type I error
bc residuals not IID
Need to use model w/ s(date)*

Generalized Additive Models Multiple Smoothers Bivariate Splines

Autocorrelation in residuals

One person's mean structure is another person's correlation structure.

If you don't have covariates to model the correlation structure, you can specify correlated errors.

This is beyond the scope of this course, but a nice tutorial is available at <https://petolau.github.io/Analyzing-double-seasonal-time-series-with-GAM-in-R/>

26/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Splines

Bivariate Splines

shifting years

30/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Splines

Bivariate Splines

Under an additive model framework, we can also consider smooth effects of two variables jointly:

$$y_i = f(x_i, z_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2).$$

We can think of $f(x_i, z_i)$ as a **surface**.

We can use 2-dimensional splines to model $f(x_i, z_i)$.

Let $\mathbf{s}_i = (x_i, z_i)$ be some pair of covariate values, and let $\mathbf{k}_m = (x_m, z_m)$ denote the m^{th} knot in the domain of x_i and z_i . We can express the smooth function as

$$f(x_i, z_i) = \beta_0 + \sum_{m=1}^M \beta_m b_m(\mathbf{s}_i; \mathbf{k}_m).$$

Note that $b_m(\cdot)$ is a basis function that maps $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$.

31/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing **Bivariate Spline**

Thin-plate Spline

One popular bivariate basis function uses **thin-plate splines**, which extends to $\mathbf{s}_i \in \mathbb{R}^d$ and $\partial^l g$ penalties. We consider $d = 2$ and $l = 2$:

$$f(\mathbf{s}_i) = \beta_0 + \beta_1 x_i + \beta_2 z_i + \sum_{m=1}^M \beta_{2+m} b_m(\mathbf{s}_i; \mathbf{k}_m)$$

using the radial basis:

$$b_m(\mathbf{s}_i; \mathbf{k}_m) = \|\mathbf{s}_i - \mathbf{k}_m\|^2 \log(\|\mathbf{s}_i - \mathbf{k}_m\|).$$

Here, $\|\mathbf{s}_i - \mathbf{k}_m\|$ is the Euclidean distance between the covariate \mathbf{s}_i and the knot location \mathbf{k}_m .

The radial basis kernel is $r^2 \log r$.

The thin-plate spline is sensitive to the scale of each variable, but invariant to rotation (isotropic).

It is best for variables measured on the same scale (e.g. geographical distance).

32/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing **Bivariate Spline**

Thin-plate Spline, cont.

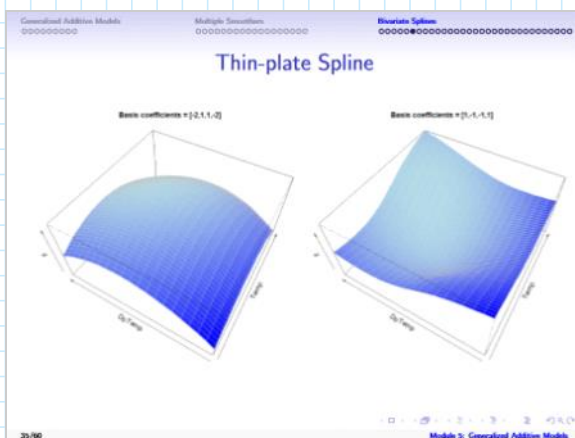
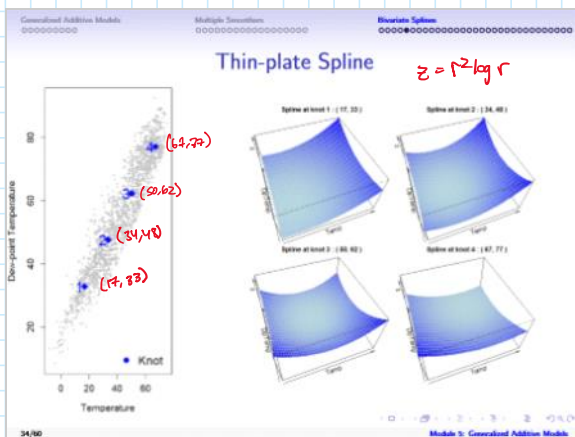
It can be shown that the thin-plate spline function minimizes

$$\sum_{i=1}^n \{y_i - f(x_i, z_i)\}^2 + \lambda \int \left(\frac{\partial^2 f(x, z)}{\partial x^2} \right)^2 + 2 \left(\frac{\partial^2 f(x, z)}{\partial x \partial z} \right)^2 + \left(\frac{\partial^2 f(x, z)}{\partial z^2} \right)^2 dx dz.$$

More information is in Wood 2017, pages 215-221, and references therein.

*Intuition: second derivatives measure wiggleness
so penalize it*

33/60 Module 5: Generalized Additive Models



Generalized Additive Models Multiple Smoothers Bivariate Splines

Thin-plate regression spline

mgcv::gam uses thin-plate regression splines as the default for smoothers of a single variable (as well as two variables).

This is implemented in a 'knot-free' manner. *NOT thin plate splines*

This is the general idea.

For $d = 2$, $l = 2$:

1. Construct the $n \times n$ matrix \mathbf{E} from $\|s_i - s_{i'}\|^2 \log(\|s_i - s_{i'}\|)$.
2. Use the singular value decomposition to find a low rank representation, e.g., k leading singular vectors, and use this in place of \mathbf{X} in the penalized objective function.
3. In practice, there are some additional things to worry about to make \mathbf{I} (for the intercept), \mathbf{x} and \mathbf{z} in the null space of the penalty.
4. Then estimate the β_0 , β_1 for \mathbf{x} , β_2 for \mathbf{z} (unpenalized) and β_3, \dots, β_k (penalized), which dramatically reduces computation costs.

The formulas for the general case (d dimensions and l th derivative) get a bit complicated; see 5.5 in Wood.

36/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Splines

Joint Effects of Temperature and Dew point Temperature

To define a bivariate smoother, simply specify $s(\text{var}_1, \text{var}_2)$ in the equation formula.

Default in mgcv::gam is to use a thin-plate regression spline.

For this demonstration, let's look at the joint effects of same-day temperature and dew point temperature on log mortality, controlling only for time trends.

specifies bivariate smooth

```
> fit1 = gam(log(alldeaths) ~ s(date2, k=100) + s(Temp, DpTemp, k = 20), data = health)
> summary(fit1)
```

Parametric coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) 4.219108 0.002963 1429 <2e-16 ***

Approximate significance of smooth terms:
edf Ref.df F p-value
s(date2) 69.527 62.144 8.501 < 2e-16 ***
s(Temp, DpTemp) 5.938 8.032 6.407 2.96e-08 ***

R-sq. (adj.) = 0.45 Deviance explained = 47.3%
GCV score = 0.016634 Scale est. = 0.015928 n = 1826

You should also use [gam.check](#) (note shown here)

37/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Splines

Joint Effects of Temperature and Dew Point Temperature

Perspective Plots from Thin-Plate Spline

38/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Splines

SKIP - direct extension of univariate splines

Tensor Product

Another way to obtain a 2-D spline is by construction. First consider the effect of a variable x_i specified by M basis functions $b_{m,x}(x_i)$

$$f(x_i) = \sum_{m=0}^M \beta_m b_{m,x}(x_i).$$

Now assume $f(x_i, z_i)$ is created by allowing each spline coefficient to vary with the second variable z_i : M

$$f(x_i, z_i) = \sum_{m=0}^M \beta_m(z_i) b_{m,x}(x_i).$$

We can express $\beta_m(z_i)$ also as a smooth function of z_i using N basis functions $b_{n,z}(z_i)$:

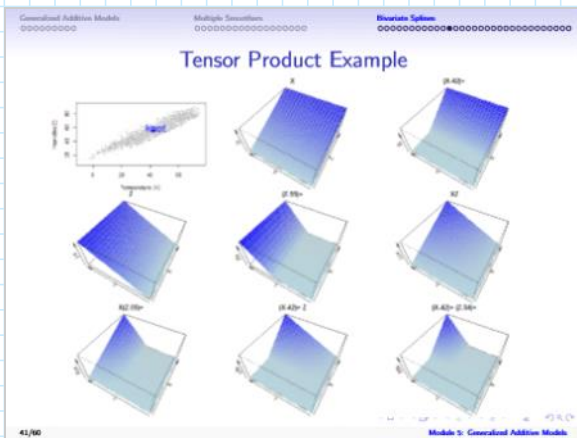
$$f(x_i, z_i) = \sum_{m=0}^M \sum_{n=0}^N \alpha_{m,n} b_{n,z}(z_i) b_{m,x}(x_i).$$

This is equivalent to expressing $f(x_i, z_i)$ as all pairwise basis functions of x_i and z_i :

$$f(x_i, z_i) = \sum_{m=0}^M \sum_{n=0}^N \alpha_{m,n} b_{m,n}(x_i, z_i).$$

39/60 Module 5: Generalized Additive Models

The slide features a purple header bar at the top. On the left, it says 'Generalized Addition Models' followed by a row of 10 small squares. In the center, it says 'Multiple Swatches' followed by a row of 10 small squares. On the right, it says 'Bivariate System' followed by a row of 10 small squares. The main title 'Tensor Product' is centered in a large, dark blue font. Below the title, the text 'Example with 1 knot: write out all pairwise interactions' is displayed. At the bottom, there is a purple footer bar containing navigation icons on the left and the text 'Module 3: Generalized Addition Models' on the right.



Multiple Smoothers
 ○○○○●○○○○○○○○○○○○○○○○

Bivariate Splines
 ○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○

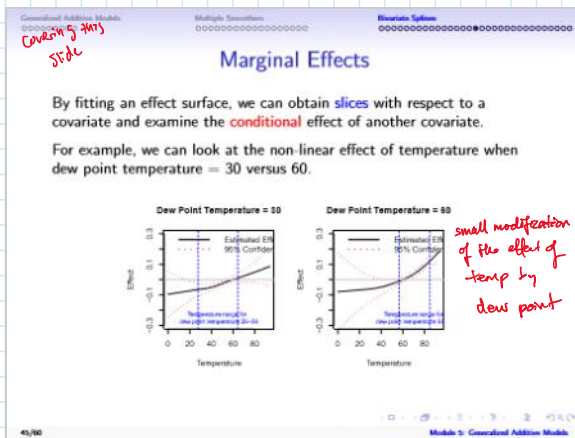
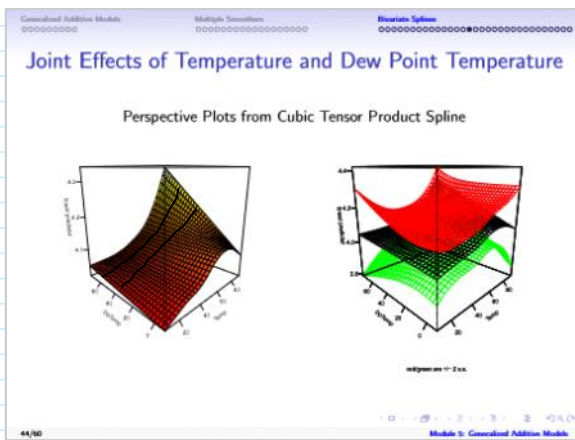
Tensor Product

The roughness penalty optimized by the tensor product is

$$\sum_{i=1}^n \{y_i - f(x_i, z_i)\}^2 + \int \lambda_x \left(\frac{\partial^2 f(x,z)}{\partial x^2} \right)^2 + \lambda_z \left(\frac{\partial^2 f(x,z)}{\partial z^2} \right)^2 dx dz .$$

- Allows penalization in each variable dimension by assigning a smoothing parameter for each marginal smooth effect.
- Choice of basis function and knots also do not need to be the same for each covariate.
- Provides a recipe for constructing flexible multivariate joint effects.
- Often used for modeling interactions where the degree of smoothness may not be the same for all covariates.

[illegible]



Generalized Additive Models Multiple Smoothing Bivariate Splines

PM_{2.5} and Temperature

Let's revisit the PM_{2.5} analysis. An important research question is whether there is an interaction between air pollution and temperature.

For this demonstration, we consider previous-day PM_{2.5} level and 3-day moving average of temperature (rmTemp).

We use thin-plate splines, which are much faster to fit than tensor; pm25.lag1 and rmTemp are very roughly on same scale.

46/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Bivariate smooth for rmTemp and PM2.5

```
> fit = gam(log(alldeaths)~s(pm25.lag1, rmTemp)+factor(dow)+s(date2, k = 75)+
+ s(lpTemp)+s(rmTemp), data = health)
> summary(fit)
```

Parametric coefficients:

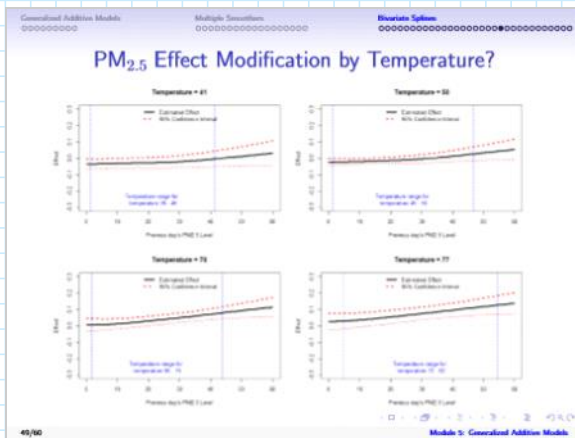
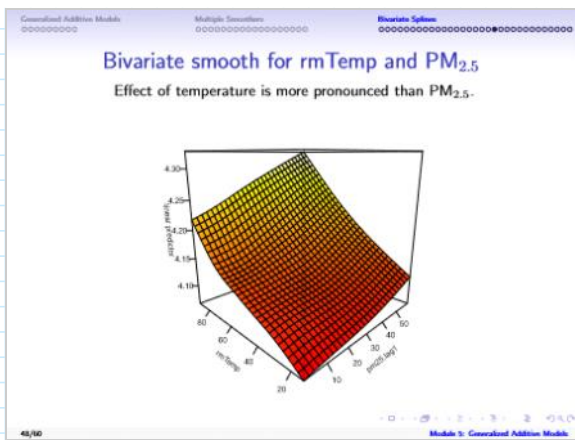
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.2191876	0.0077934	541.380	<2e-16 ***
factor(dow)Friday	-0.0142658	0.0110132	-1.296	0.1964
factor(dow)Monday	0.0222996	0.0110146	2.025	0.0431 *
factor(dow)Saturday	-0.0092121	0.0110082	-0.837	0.4028
factor(dow)Thursday	0.0011288	0.0110248	0.102	0.9185

Approximate significance of smooth terms:

	edf	Ref.df	F	p-value
s(pm25.lag1,rmTemp)	5.858	8.381	4.048	6.35e-05 ***
s(date2)	57.808	66.508	6.721	< 2e-16 ***
s(lpTemp)	1.362	1.633	35.075	8.11e-10 ***
s(rmTemp)	1.000	1.000	35.085	3.79e-09 ***

R-sq. (adj.) = 0.456 Deviance explained = 47.8%
 GCV score = 0.016379 Scale est. = 0.015724 n = 1825

47/60 Module 5: Generalized Additive Models



Generalized Additive Models Multiple Smoothers Bivariate Spline

Extracting Effects from Smooth Function

- Let's say we are interested in the difference in $f(x_i)$ when x_i increases from value a to value b . $\rightarrow (40, 70) \text{ to } (50, 70)$

$$\hat{f}(b) - \hat{f}(a) = \sum_{m=1}^M \hat{\beta}_m b_m(b) - \sum_{m=1}^M \hat{\beta}_m b_m(a) = \sum_{m=1}^M \hat{\beta}_m \{b_m(b) - b_m(a)\}.$$

The covariance matrix is given by $Cov(\mathbf{B}_{b-a}\hat{\beta}) = \mathbf{B}_{b-a}' Cov(\hat{\beta}) \mathbf{B}_{b-a}$, where

$$\mathbf{B}_{b-a} = [b_1(b) - b_1(a), b_2(b) - b_2(a), \dots, b_M(b) - b_M(a)]'$$

- However $b_m(\cdot)$'s are basis functions and $b_m(b) - b_m(a) \neq b_m(b - a)$. We will need to construct the new covariate vector \mathbf{B} ourselves.

50/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothers Bivariate Spline

Extracting Effects from a Smooth Function

We want to estimate the effect of a 10 $\mu\text{g}/\text{m}^3$ increase in PM_{2.5} levels from 40 to 50 $\mu\text{g}/\text{m}^3$ when temperature = 70F.

```
> X1 = predict(fit, data.frame(pm25.lag1 = 40, rmTemp=70,
+   rmTemp=0, rmUpTemp = 0, dow = "Sunday", date2=0), type = "lmatrix")
> X2 = predict(fit, data.frame(pm25.lag1 = 50, rmTemp=70,
+   rmTemp=0, rmUpTemp = 0, dow = "Sunday", date2=0), type = "lmatrix")
```

B₀ from previous slide along w/ other variables

```
> X.diff = X2 - X1
> dim(X.diff)
[1] 1 128
> Est = X.diff %*% coef(fit) ## Estimate
> se = sqrt(X.diff %*% vcov(fit) %*% t(X.diff)) ## Standard Error
> Est; se
[1]
1 0.02213608
1
1 0.01096892
```

the whole design matrix

formula for combining variances

So our estimate is 0.022 (95%CI 0.001, 0.044).

The same effect of a 10-unit increase in PM_{2.5} levels from 20 to 30 $\mu\text{g}/\text{m}^3$ is 0.020 (95%CI 0.004, 0.036).

These two estimates are very similar because the PM_{2.5} effect appears to be quite linear.

51/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Inference in models with smoothing splines

For parametric terms, the inference is identical to purely parametric model.

You can use the t-statistics from `summary(fit)`.

```
> fit.full ~ gam(log(cr6plus) ~ s(pm25.lag1, rmTemp) + fdow + s(date2, k = 100) + s(DpTemp) + s(rndpTemp))
> summary(fit.full)
```

Family: gaussian
Link function: identity

Formula:
log(cr6plus) ~ s(pm25.lag1, rmTemp) + fdow + s(date2, k = 100) + s(DpTemp) + s(rndpTemp)

Parametric coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.219093	0.007784	541.997	<2e-16 ***
fdowFriday	-0.014166	0.011001	-1.288	0.1980
fdowMonday	0.022352	0.011002	2.032	0.0423 *
fdowSaturday	-0.009111	0.010996	-0.829	0.4074
fdowThursday	0.001261	0.011012	0.114	0.9089
fdowTuesday	-0.001689	0.011002	-0.153	0.8780
fdowWednesday	0.000470	0.011012	0.043	0.9660

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

*nothing changes over all the usual SEs, no change from 2ms
can test significance of ?*

Generalized Additive Models Multiple Smoothing Bivariate Splines

Inference in models with smoothing splines

For overall effect of fdow, use anova:

```
> anova(fit.full)
```

Family: gaussian
Link function: identity

Formula:
log(cr6plus) ~ s(pm25.lag1, rmTemp) + fdow + s(date2, k = 100) + s(DpTemp) + s(rndpTemp)

Parametric Terms:

	df	F	p-value
fdow	6	2.154	0.0448

Generalized Additive Models Multiple Smoothing Bivariate Splines

Inference in models with smoothing splines

For smoothed terms, the inference is approximate because the distribution of the test statistics is impacted by the penalization.

The idea is to jointly test the significance of the β_j for the coefficients corresponding to the spline of the j th smooth term.

There is some discussion of these approximate p-values in `?summary.gam`. For detailed discussion, see Wood 2017 p.304.

Approximate inference for smooth effects:

Approximate significance of smooth terms:

	edf	Ref.df	F	p-value
s(pm25.lag1, rmTemp)	5.433	7.738	4.246	6.30e-06 ***
s(date2)	62.047	74.627	6.118	< 2e-16 ***
s(DpTemp)	1.210	1.379	29.260	2.89e-09 ***
s(rndpTemp)	1.000	1.000	34.655	4.68e-09 ***

R-sq.(adj) = 0.458 Deviance explained = 45%
CV = 0.016377 Scale est. = 0.015688 n = 1826

Generalized Additive Models Multiple Smoothing Bivariate Splines

Inference in models with smoothing splines

Should we use a bivariate spline for pm25.lag1 and rmTemp, or two univariate splines? Don't do this:

```
> anova(fit.reduced2, fit.full, test="F")
```

Analysis of Deviance Table

Model 1: log(cr6plus) ~ s(pm25.lag1) + s(rmTemp) + s(date2, k = 100) + fdow + s(DpTemp) + s(rndpTemp)
Model 2: log(cr6plus) ~ s(pm25.lag1, rmTemp) + fdow + s(date2, k = 100) + s(DpTemp) + s(rndpTemp)

	Resid. Df	Resid. Dev	Df	Deviance	F	Pr(>F)
1	1734	9095	27	353616		
2	1733	2552	27	428321	1.6542684	-0.074706162

What happened here? →

Models are not nested.

Using tensor splines, we can construct in a special way to make nested. Then we can recast the problem as testing for an interaction.

*note: this is univariate splines for pm25.lag1 and rmTemp
note: this is different we had interaction between continuous and categorical variable →*

Generalized Additive Models Multiple Smoothing Bivariate Splines

Testing for an interaction

To test for an interaction between `pm25.lag1` and `rmTemp`, we can construct the tensor spline in a special way such that the univariate splines are in the null space of the penalty of the tensor spline. See pages 243 and the example on page 343-346.

```
fit.full.tensor = gam(log(cr66plus) ~ s(pm25.lag1, bs="cr") + s(rnTemp, bs="cr") +
  ti(pm25.lag1, rnTemp, bs="cr") + fdow + s(date2, k = 100) + s(DpTemp) + s(rnDpTemp), data = health)
summary(fit.full.tensor)
anova(fit.full.tensor)
```

16/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

Testing for an interaction

```
> anova(fit.full.tensor)
```

Family: gaussian
Link function: identity

Formula:
`log(cr66plus) ~ s(pm25.lag1, bs = "cr") + s(rnTemp, bs = "cr") +`
`ti(pm25.lag1, rnTemp, bs = "cr") + fdow + s(date2, k = 100) +`
`s(DpTemp) + s(rnDpTemp)`

Parametric Terms:

	df	F	p-value
fdow	6	2.24	0.0371

Approximate significance of smooth terms:

	edf	Ref. df	F	p-value
<code>s(pm25.lag1)</code>	1.000	1.001	4.147	0.0418
<code>s(rnTemp)</code>	5.301	6.480	2.524	0.0171
<code>ti(pm25.lag1, rnTemp)</code>	1.000	1.000	3.463	0.0629
<code>s(date2)</code>	59.917	72.436	6.143	<2e-16
<code>s(DpTemp)</code>	1.000	1.000	41.860	<2e-16
<code>s(rnDpTemp)</code>	1.000	1.000	31.651	<2e-16

$H_0: f_{x,z}(x_i, z_i)$ does not improve model fit over $f_x(x_i)$ and $f_z(z_i)$.
p=0.06. *> 0.05, model does not contribute to fit*

17/60 Module 5: Generalized Additive Models

CHECK THIS
LOOK FOR
WEIRD
SYNTAX

Generalized Additive Models Multiple Smoothing Bivariate Splines

Inference in models with smoothing splines

Let's refit the model with `pm25.lag1` as a linear term.

```
> fit.reduced3 = gam(log(cr66plus) ~ pm25.lag1 + s(rnTemp) +
  s(date2, k = 100) + fdow + s(DpTemp) + s(rnDpTemp), data = health)
> summary(fit.reduced3)
```

Family: gaussian
Link function: identity

Formula:
`log(cr66plus) ~ pm25.lag1 + s(rnTemp) + s(date2, k = 100) + fdow +`
`s(DpTemp) + s(rnDpTemp)`

Parametric coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.203e+00	1.018e-02	412.689	<2e-16 ***
<code>pm25.lag1</code>	1.073e-03	4.508e-04	2.381	0.0174 *
<code>fdowFriday</code>	-1.478e-02	1.098e-02	-1.345	0.1789
<code>fdowMonday</code>	2.248e-02	1.098e-02	2.048	0.0407 *
<code>fdowSaturday</code>	-9.026e-03	1.098e-02	-0.822	0.4110
<code>fdowThursday</code>	5.425e-04	1.100e-02	0.049	0.9607
<code>fdowTuesday</code>	-1.433e-03	1.098e-02	-0.130	0.8962
<code>fdowWednesday</code>	9.336e-05	1.100e-02	0.008	0.9932
...				

R-sq. (adj.) = 0.459 Deviance explained = 48.2%
GCV = 0.016324 Scale est. = 0.015642 n = 1825

18/60 Module 5: Generalized Additive Models

Generalized Additive Models Multiple Smoothing Bivariate Splines

GAM or linear?

It may seem like we can test for whether or not to include a smooth term versus linear term using anova. However, the anova test can produce funny results when the EDF is approximately one, as the change in DF is very small. So, I suggest not doing this:

```
> anova(fit.reduced2, fit.reduced3, test="F")
```

Analysis of Deviance Table

	Model 1: <code>log(cr66plus) ~ s(pm25.lag1) + s(rnTemp) + s(date2, k = 100) +</code> <code>fdow + s(DpTemp) + s(rnDpTemp)</code>	Model 2: <code>log(cr66plus) ~ pm25.lag1 + s(rnTemp) + s(date2, k = 100) + fdow +</code> <code>s(DpTemp) + s(rnDpTemp)</code>				
	Resid. Df	Resid. Dev	Df	Deviance	F	Pr(>F)
1	1734.9	27.354				
2	1734.9	27.353	0.016157	0.00029882	1.1824	0.0325 *

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

19/60 Module 5: Generalized Additive Models

Generalized Additive Models
○○○○○○○○○

Multiple Smoothing
○○○○○○○○○○○○○○○○○○○

Bivariate Splines
○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○●

GAM or linear?

When you have an effect that is close to linear, your interpretation is very similar whether or not you use a gam, and when $EDF=1$, then the approximate p values are equivalent to refitting with a linear term.

To determine whether or not to include a smooth or linear effect, I suggest looking at the EDF. If it is greater than 1, than it seems reasonable to use a smooth term.

There is not really any disadvantage to modeling with a smooth term, except for some extra work we need to do for interpretation.

One approach to test for a non-linear effect is to construct a special spline that separates the linear and non-linear parts:
<https://stats.stackexchange.com/questions/449641/is-there-a-hypothesis-test-that-tells-us-whether-we-should-use>

06/60

Module 5: Generalized Additive Models