

Module 5 Part 1: Splines

Wednesday, October 25, 2023 14:29



BIOS526_M
5_PartI_S...

Module 5, Part 1: Basis Expansion and Parametric Splines, Bias-Variance Trade-off, and Model Selection

BIOS 526

1/40 M5 - Part 1

Reading

- Chapter 5 (Sections 5.1 and 5.2 on parametric splines) of Hastie et al. [Elements of Statistical Learning](#).
- Chapter 5 of James et al. (Cross-validation) [Introduction to Statistical Learning](#)

Concepts

- Basis functions and knots.
- Piecewise linear splines.
- Piecewise cubic splines.
- **Bias-Variance Tradeoff.**
- Criteria for model prediction performance (CV, GCV, AIC).

2/40 M5 - Part 1

Parametric Splines

Model Selection

Parametric Splines

3/40

MS - Part 1

Parametric Splines

Model Selection

Motivating Example: Birth Weight and Mother's Age

- gw: gestational age in weeks
- age: maternal age at delivery.
- bw: birth weight at delivery in grams.

```
> load ("GAbirth.RData")
> str (dat)
'data.frame': 5000 obs. of 3 variables:
 $ gw : num 41 38 39 40 40 38 40 38 37 ...
 $ age: int 33 23 26 30 26 18 19 37 23 25 ...
 $ bw : num 3540 3287 3438 3419 3278 ...
```

4/40

MS - Part 1

Parametric Splines

Model Selection

(birth weight → gest. week)

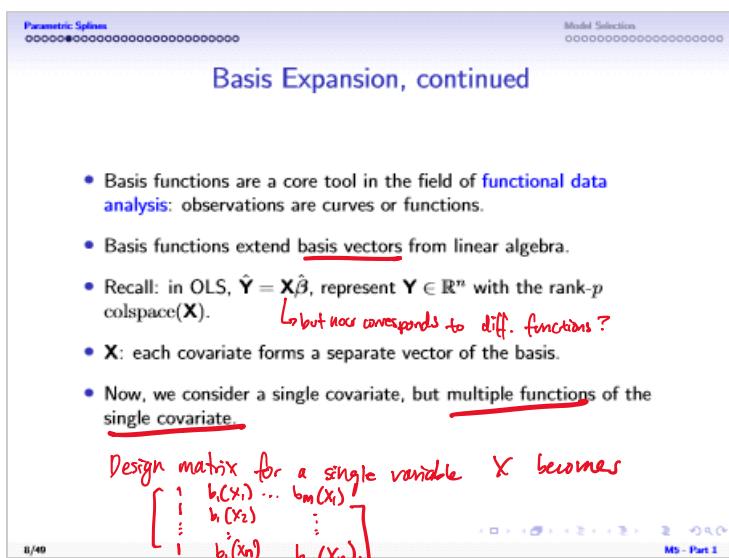
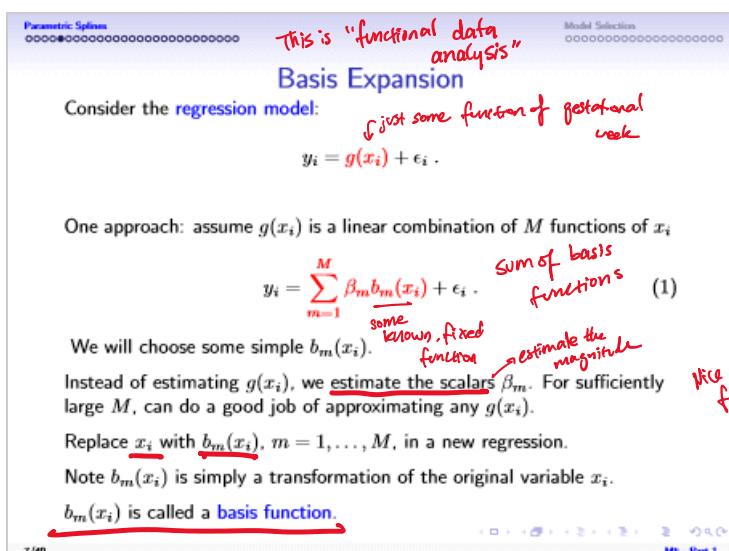
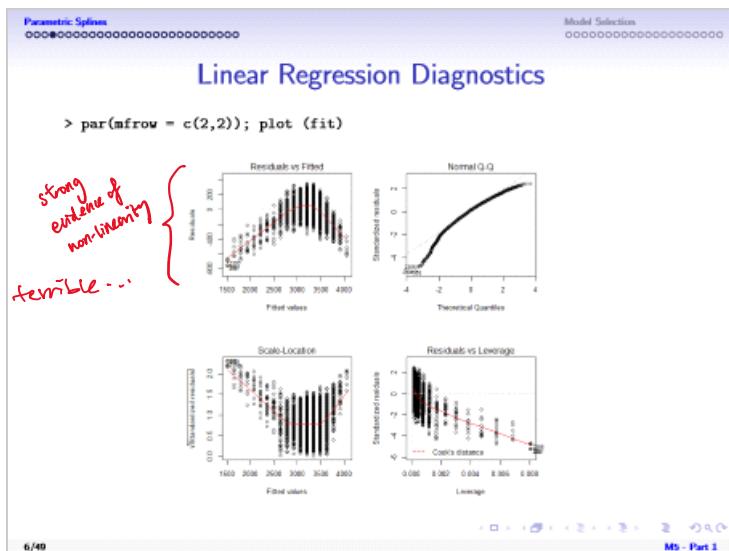
Linear Regression

```
> fit = lm (bw~gw, data = dat)
> summary(fit)
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -2177.513    39.361   -55.34 <2e-16 ***
gw          141.808     1.019   139.10 <2e-16 ***

> plot (dat$bw~dat$gw, xlab = "Gestational Week", ylab = "Birth weight")
> abline(fit, col="red", lwd = 2)
```

5/40

MS - Part 1



Concrete example

Basis Expansion Examples

Polynomial

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \varepsilon_i$$

$$g(x) = \sum \beta_n b_n(x) ??$$

$$b_0(x_i) = 1 \quad b_1(x_i) = x_i \quad b_2(x_i) = x_i^2 \quad b_3(x_i) = x_i^3$$

Indicator

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 I_{x_i > 0} + \varepsilon_i$$

jump in predicted \hat{y} 's

$$b_0(x_i) = 1 \quad b_1(x_i) = x_i \quad b_2(x_i) = I_{x_i > 0} \begin{cases} 1 & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

Periodic

$$y_i = \beta_0 + \beta_1 \sin(x_i/K) + \beta_2 \cos(x_i/K) + \varepsilon_i$$

$$b_0(x_i) = 1 \quad b_1(x_i) = \sin(x_i/K) \quad b_2(x_i) = \cos(x_i/K)$$

MS - Part 1

Polynomial Regression

creates x^2

```
> fit2 = lm(bw~gw+I(gw^2), data = dat) Quadratic polynomial
> fit3 = lm(bw~gw+I(gw^2) + I(gw^3), data = dat) Cubic polynomial
```

BAD TAIL BEHAVIOR

Note that polynomial functions often cannot model the ends very well.

MS - Part 1

Piecewise Regression aka step functions

Growth rates may differ

To model non-linear relationship: divide range of the covariate into some suitable regions \rightarrow these will define our knots

E.g., pregnancy length: preterm (< 37 weeks), full-term (37-41 weeks), and post-term (> 42 weeks).

Model with polynomial functions separately within each region. The interior values that define the regions are called **knots**.

MS - Part 1

Parametric Splines

Model Selection

Piecewise Regression Specification

Piecewise regression is equivalent to a model where the basis functions of x_i interact with dummy variables for regions.

E.g., Piecewise quadratic regression:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 D_{1i} + \beta_4 x_i \times D_{1i} + \beta_5 x_i^2 \times D_{1i} + \beta_6 D_{2i} + \beta_7 x_i \times D_{2i} + \beta_8 x_i^2 \times D_{2i}$$

where

$$D_{1i} = \begin{cases} 1 & 37 \leq x_i < 42 \\ 0 & \text{otherwise} \end{cases} \quad D_{2i} = \begin{cases} 1 & x_i \geq 42 \\ 0 & \text{otherwise} \end{cases}$$

We only need two dummy variables for three regions.

MS - Part 1

Parametric Splines

Model Selection

Piecewise Regression

- The positive association between pregnancy length and birth weight decreases for higher gestational weeks.
- Here, linear, quadratic, and cubic result in similar trend in each region.

Birth weight

Gestational Week

MS - Part 1

Parametric Splines

Model Selection

Piecewise Splines

Piecewise polynomials have limitations:

- The regression functions do not match at knot locations.
- Requires many regression coefficients (degrees of freedom).

Splines are basis functions for piecewise regressions that are connected at the interior knots.
i.e. will be continuous

Splines are continuous, which makes them aesthetically appealing.

Birth weight

MS - Part 1

Linear Splines

A linear piecewise spline at knot locations 36.5 and 41.5 is specified as
 $y_i = \beta_0 + \beta_1 x_i + \beta_2(x_i - 36.5)_+ + \beta_3(x_i - 41.5)_+$
 one intercept
 6 "global" slope

Here we only need 4 regression coefficients, instead of 6 in a model that does not force the regression lines to connect at knots.

Count the parameters: 2 parameters per line * 3 lines - 2 constraints = 4 where $(\dots)_+$ denotes the positive part.

$$(x_i - 36.5)_+ = \begin{cases} x_i - 36.5 & \text{if } x_i \geq 36.5 \\ 0 & \text{otherwise} \end{cases}$$

$$(x_i - 41.5)_+ = \begin{cases} x_i - 41.5 & \text{if } x_i \geq 41.5 \\ 0 & \text{otherwise} \end{cases}$$

MS - Part 1

Linear Splines: Coefficient Interpretation

$y_i = \beta_0 + \beta_1 x_i + \beta_2(x_i - 36.5)_+ + \beta_3(x_i - 41.5)_+$

For $x_i < 36.5$:
 $y_i = \beta_0 + \beta_1 x_i$.

For $36.5 \leq x_i < 41.5$:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \beta_2(x_i - 36.5) \\ &= (\beta_0 - \beta_2 * 36.5) + (\beta_1 + \beta_2)x_i. \end{aligned}$$

For $41.5 \geq x_i$:
 changes intercept
 (jumps together \hat{y}_i to make it continuous)
 changes slope

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \beta_2(x_i - 36.5) + \beta_3(x_i - 41.5) \\ &= (\beta_0 - \beta_2 * 36.5 - \beta_3 * 41.5) + (\beta_1 + \beta_2 + \beta_3)x_i. \end{aligned}$$

β_2 and β_3 represent changes in slope compared to the previous region

MS - Part 1

See MG_part1_Splines.R

Linear Splines: Coefficient Interpretation

```
> Sp1 = (dat$gw - 36.5)*as.numeric(dat$gw >= 36.5)
> Sp2 = (dat$gw - 41.5)*as.numeric(dat$gw >= 41.5)
> fit7 = lm (bw ~ gw+Sp1 + Sp2, data = dat)
> summary (fit7)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-5013.768	62.747	-79.90	<2e-16 ***
gw	222.620	1.757	126.71	<2e-16 ***
Sp1	-121.427	2.549	-47.64	<2e-16 ***
Sp2	-152.948	10.492	-14.58	<2e-16 ***

- A 223 g (CI_{95%} 219, 226) increase in birth weight was associated with a one week increase in pregnancy length prior to week 36.5. (gestation time)
- The rate of increase then drops by 121.4 g/week to 101.2 g/week (CI_{95%} 99, 104) between pregnancy weeks 36.5 to 41.5. (see R code)
- After pregnancy week 41.5, birth weight was negatively associated with gestational week by -52 g/week (CI_{95%} -72, -32).

slope $\pm 1.96(\text{SE})$

MS - Part 1

Cubic Splines

Linear splines are easy to interpret → e.g. like a shift in growth rate
 The resulting function $g(x_i)$ is not smooth. The first-derivative of BUT

MS - Part 1

$$\text{Cor}(\hat{\beta}) = C' \text{Var}(\hat{\beta}) C$$

where \leftarrow

selection region?
 A vector of constants
 \downarrow
 e.g.
 $c(0, 1, 1, 0)$

↓
 variance covariance matrix of the coefficients

Linear splines are easy to interpret. → e.g. like a shift in growth rate
 BUT
 The resulting function $g(x_i)$ is not smooth. The first-derivative of

$$g(x_i) = \beta_0 + \beta_1 x_i + \beta_2 (x_i - 36.5)_+ + \beta_3 (x_i - 41.5)_+$$

 is discontinuous at the knots.

We can also work with piecewise cubic functions that are connected at the knots:

$$g(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 (x_i - 36.5)_+^3 + \beta_5 (x_i - 41.5)_+^3$$

Cubic splines are popular because:

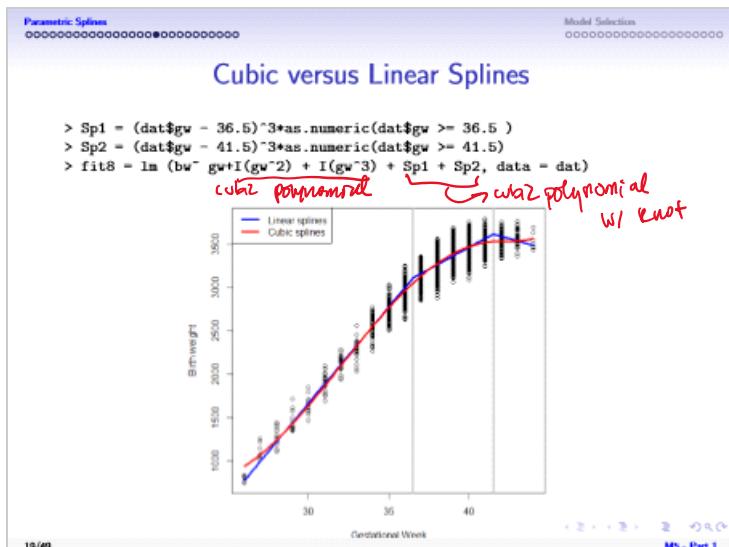
- continuous 2nd-derivatives, typically what we view as smooth visually.

"Folks like these bc they look nice!"

18/49

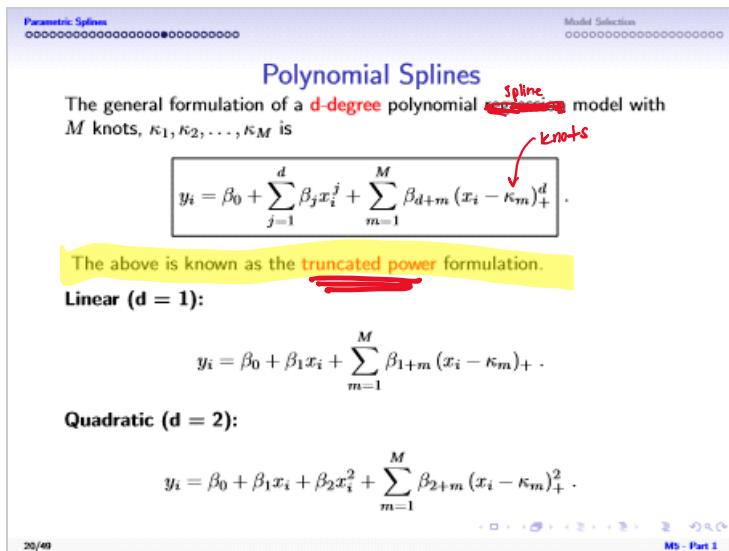
c.g.
 $c(0, 1, 1, 0)$
 will select
 β_1 & β_2

MS - Part 1



19/49

MS - Part 1



20/49

MS - Part 1

Parametric Splines Model Selection

Interpreting Polynomial Splines

- Linear splines: a good choice for interpretability – the coefficients at the knots represent the change in slope from the previous region.
- Cubic splines: appears smooth, which is often biologically reasonable. [Generally we do not interpret the individual coefficients of the truncated polynomials. **for d>1**]

```
lm(formula = bw ~ gw + I(gw^2) + I(gw^3) + Sp1 + Sp2, data = dat)
  Estimate Std. Error t value Pr(>|t|)
(Intercept) 2.890e+04 3.717e+03 7.776 9.02e-15 ***
gw          -2.968e+03 3.335e+02 -8.901 < 2e-16 ***
I(gw^2)      9.972e+01 9.900e+00 10.073 < 2e-16 ***
I(gw^3)     -1.036e+00 9.734e-02 -10.641 < 2e-16 ***
Sp1         7.732e-01 2.482e-01 3.115 0.00185 **
Sp2         7.455e+00 3.471e+00 2.148 0.03177 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 109.8 on 4994 degrees of freedom
 Multiple R-squared: 0.881, Adjusted R-squared: 0.8809
 F-statistic: 7397 on 5 and 4994 DF, p-value: < 2.2e-16

21/49 MS - Part 1

Parametric Splines Model Selection

B-Splines \rightarrow lots of "incremental improvements"

The truncated power formulation can sometimes experience numerical problems when fitting because:

- the covariates are highly correlated;
- d^{th} power can be small or large with large d .

B-splines can be used to represent the same space as truncated polynomial splines. Just a diff. way of representing polynomial splines.
You'll get the same fit

Mathematically, these will produce the same predictions as truncated polynomial splines.

Computationally, may be more accurate – avoid overflow errors.

Scaled between 0 and 1. Provide better numerical stability.

Commonly implemented in statistical software but in my experience you can just use cubic splines with modern computers.

22/49 MS - Part 1

Parametric Splines Model Selection

B-Splines and computation

The `bs()` function in R will create the appropriate design matrix.

```
### Load the splines package
library(splines)

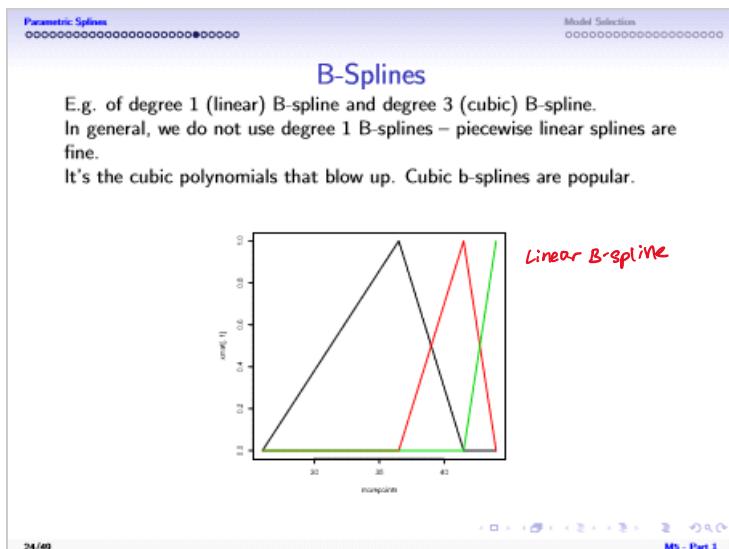
### Linear splines
fit1 = lm ( bw ~ bs(gw, knots = c(36.5, 41.5), degree = 1), data = dat)

### Quadratic splines
fit2 = lm ( bw ~ bs(gw, knots = c(36.5, 41.5), degree = 2), data = dat)

### Cubic splines
fit3 = lm ( bw ~ bs(gw, knots = c(36.5, 41.5), degree = 3), data = dat)
```

linear quadratic cubic

23/49 MS - Part 1



Parametric Splines Model Selection

B-Splines

Coefficients are not interpretable.

```
lm(formula = bw ~ bs(gw, knots = c(36.5, 41.5), degree = 3),
   data = dat)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-399.13	-74.61	4.47	77.06	301.40

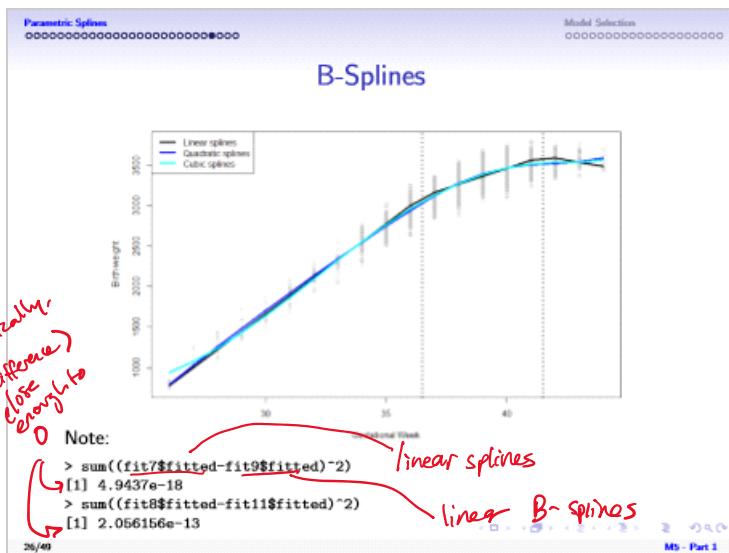
Coefficients: *Cubic version*

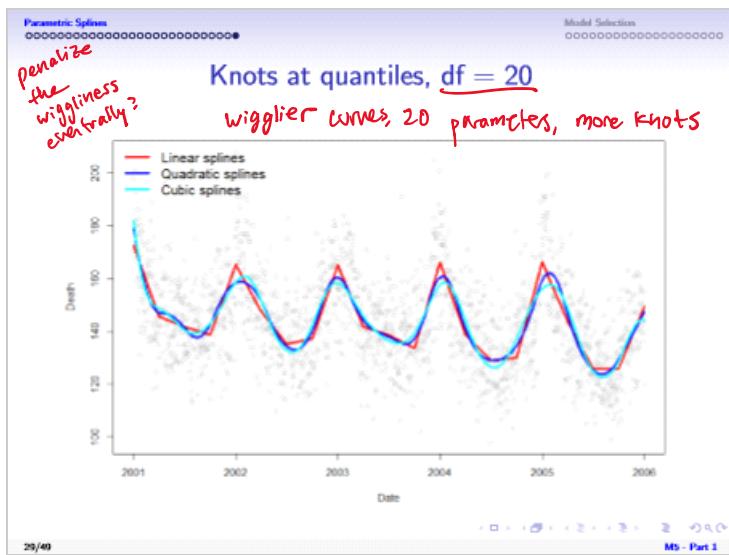
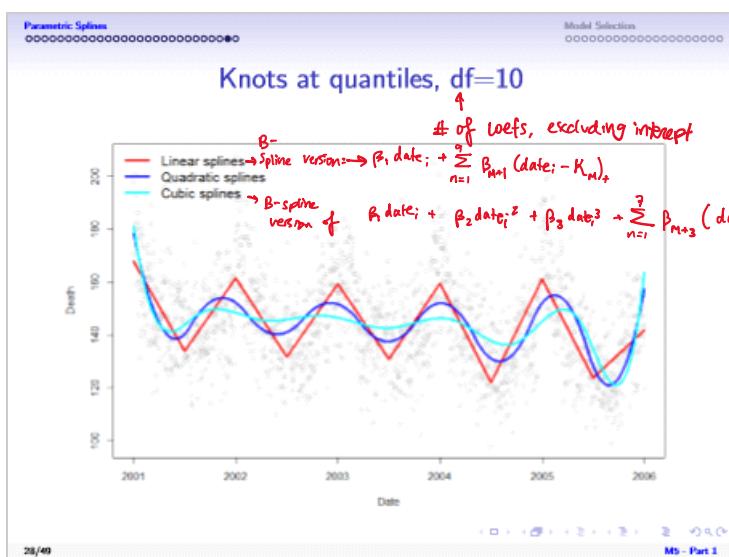
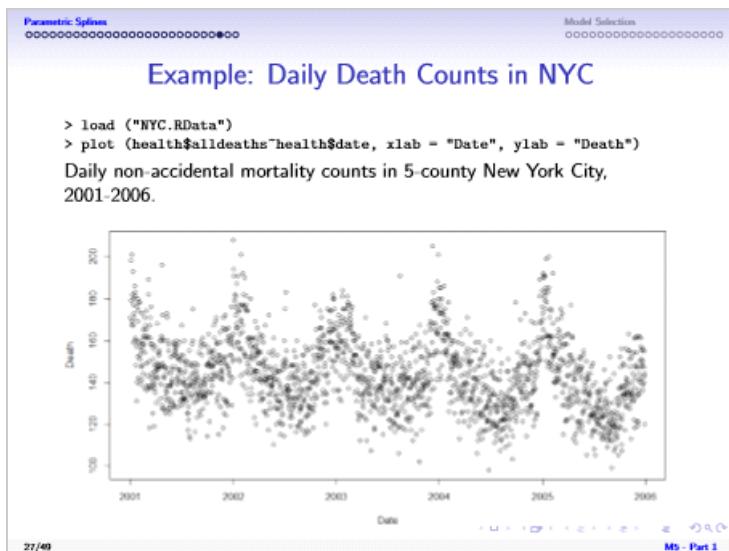
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	937.45	34.65	27.056	< 2e-16
bs(gw, knots = c(36.5, 41.5), degree = 3)1	408.19	57.92	7.048	2.07e-12
bs(gw, knots = c(36.5, 41.5), degree = 3)2	2037.52	32.46	62.779	< 2e-16
bs(gw, knots = c(36.5, 41.5), degree = 3)3	2655.45	38.15	69.608	< 2e-16
bs(gw, knots = c(36.5, 41.5), degree = 3)4	2582.17	35.94	71.851	< 2e-16
bs(gw, knots = c(36.5, 41.5), degree = 3)5	2633.37	52.76	49.912	< 2e-16

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

Residual standard error: 109.8 on 4994 degrees of freedom
 Multiple R-squared: 0.881, Adjusted R-squared: 0.8809
 F-statistic: 7397 on 5 and 4994 DF, p-value: < 2.2e-16

25/49 MS - Part 1





Parametric Splines

Model Selection

Model Selection

30/49

Ms - Part 1

Parametric Splines

Model Selection

Model Selection

We need to pick:

1. Which basis function? Linear versus cubic. → don't want too much, well use NCV
2. How many knots?
3. Where to put the knots?

Challenge: these models are usually **not nested!**

Main Ideas:

- Choice of basis function usually does not have a large impact on model fit, especially when there are enough knots and $g(x)$ is smooth.
- When there are enough knots, their locations are less important too.

For now, we'll focus on the **how many** question.

31/49

Ms - Part 1

Parametric Splines

Model Selection

Bias-Variance Trade-Off

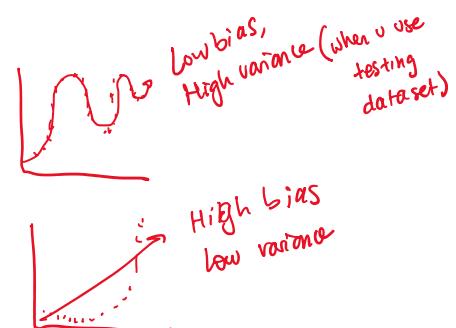
More knots → { Can better capture fine-scale trends
Need to estimate more coefficients → so we risk overfitting

Cubic Splines with Different Number of Knots

32/49

Ms - Part 1

what model to use?
Cross Validation can help us pick.



Parametric Splines Model Selection

Towards the Bias-Variance Decomposition

Assume Y comes from some true model

$$Y = g(X) + \epsilon, \quad \epsilon \sim (0, \sigma^2).$$

start w/ non-parametric regression
where $g(\cdot)$ is some unknown function.

For given x_i , we estimate Y with $\hat{g}(x_i)$, where here we assume $\hat{g}(x_i) \perp\!\!\!\perp \epsilon$. In other words, x_i is an "unseen" data point.

The expected squared difference between Y and the estimator $\hat{g}(X)$ is the population mean squared error (MSE). For given x_i ,

population MSE

$$\begin{aligned} E[\{Y - \hat{g}(X)\}^2 | X = x_i] &= \text{derive...} \\ &= \sigma^2 + E[\{g(x_i) - \hat{g}(x_i)\}^2] \\ &\quad \downarrow \quad \text{true mean} \quad \text{estimate} \\ &\quad \text{don't know} \end{aligned}$$

33/49 MS - Part 1

Parametric Splines Model Selection

Derivations

$$\begin{aligned} E[Y - \hat{g}(x)]^2 &= E[(Y - g(x_i) + g(x_i) - \hat{g}(x_i))^2] \\ &= E[(Y - g(x_i))^2] + 2E(Y - g(x_i))(g(x_i) - \hat{g}(x_i)) + E[g(x_i) - \hat{g}(x_i)]^2 \\ &\quad \downarrow \quad \text{true mean} \\ &= \sigma^2 + 2E[\epsilon_i(g(x_i) - \hat{g}(x_i))] + E[g(x_i) - \hat{g}(x_i)]^2 \\ &= \sigma^2 + 2E[\epsilon_i] E[g(x_i) - \hat{g}(x_i)] + E[g(x_i) - \hat{g}(x_i)]^2 \\ &= \boxed{\sigma^2 + E[g(x_i) - \hat{g}(x_i)]^2} \\ &\quad \text{can't do anything about } \sigma^2 \end{aligned}$$

34/49 MS - Part 1

Parametric Splines Model Selection

Bias-Variance Decomposition, cont.

Note: \hat{g} is estimated from g , which is random

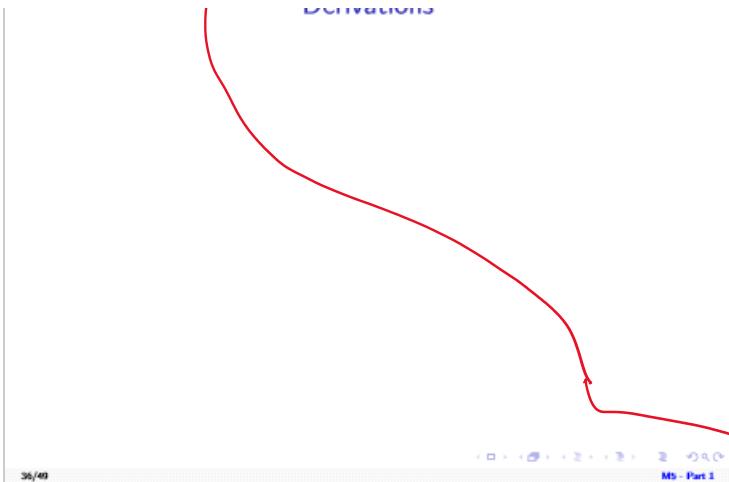
$$\begin{aligned} E\{g(x_i) - \hat{g}(x_i)\}^2 &= \dots \\ &= \underbrace{E[g(x_i)^2]}_{\downarrow} - 2E[g(x_i)\hat{g}(x_i)] + E[\hat{g}(x_i)^2] \\ &= g(x_i)^2 - 2g(x_i)E\hat{g}(x_i) + E[\hat{g}(x_i)^2] \star \\ \text{Bias}^2 &= [g(x_i) - E\hat{g}(x_i)]^2 = g(x_i)^2 - 2g(x_i)E\hat{g}(x_i) + [E\hat{g}(x_i)]^2 \\ \text{Variance} &= E[\hat{g}(x_i)^2] - E[g(x_i)]^2 \\ \text{Bias}^2 + \text{Variance} &= g(x_i)^2 - 2g(x_i)E\hat{g}(x_i) + [E\hat{g}(x_i)]^2 + E[\hat{g}(x_i)^2] - E[g(x_i)]^2 \\ &= g(x_i)^2 - 2g(x_i)E\hat{g}(x_i) + E[\hat{g}(x_i)^2] \\ &= \star \\ &= \boxed{\text{MSE}} \end{aligned}$$

35/49 MS - Part 1

Parametric Splines Model Selection

Derivations

Derivations



Parametric Splines

Model Selection

Cross-validation Error

(you also think $\hat{g}^{(-i)}(x_i)$: estimate g using all data except x_i , then evaluate at x_i)

Let $\hat{g}_i^{(-i)}$ be the fitted value of y_i at x_i using all the data except y_i . The ordinary leave-one-out cross-validation (LOOCV) is given by

$$CV = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{g}_i^{(-i)})^2$$

Intuitively, CV estimates the population MSE, $E[(Y - \hat{g}(X))^2]$, with a sample MSE.

Note here we are averaging across the values of x_i .

Caveat from Hastie et al: "Discussions of error rate estimation can be confusing, because we have to make clear which quantities are fixed and which are random..."

Expectation w.r.t respect to dist. of Y ?

Treat X as fixed (e.g. X is defined by a grid of points, like dates)

$$E[E[(Y - g(x))^2 | X]]$$

Parametric Splines

Model Selection

Overfitting

- If you **overfit**, then you start to fit the noise in the data, i.e., you estimate $g(x_i) + \epsilon_i$, instead of $g(x_i)$.
- With overfitting, new data are poorly predicted.
- From Bias-Variance perspective: overfitting = small bias, high variance $\rightarrow E[\hat{g}(x) - E[g(x)]^2 \text{ is large}$ $\rightarrow [g(x_i) - E[g(x_i)]^2 \text{ is small}$ *unseen date point y_i will be predicted poorly*
- In splines, lines that are very wiggly = high variance.
- If you **underfit**, then you tend to estimate a smoothed version of $g(x_i)$, at extreme, fit a mean s.t. $\hat{g}(x_i) = \frac{1}{n} \sum y_i \rightarrow$ then, $(g(x_i) - \hat{y}_i)^2$ will be large.
- From Bias-Variance perspective: underfitting = large bias, small variance $E(\frac{1}{n} \sum y_i - E[\frac{1}{n} y_i])^2$ will be small
- In splines, lines that vary little = low variance.

Parametric Splines

Model Selection

Generalized Cross-validation Error

One variation of LOOCV is the generalized cross-validation (GCV)

*generalized
for
penalized regression's*

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

Compared to CV,

- GCV replaces each H_{ii} by the average of all H_{ii} . i.e. $\frac{1}{n} \sum_{i=1}^n H_{ii}$
- GCV is a weighted version of CV.
- For complicated models (not OLS), CV is computationally costly, as it requires fitting the model many times.
- GCV often easier to calculate. can always do $\hat{Y} = SY$

We want to pick a model which minimizes CV or GCV.

Note: k-fold
not part of exam
but GCV will be

K-Fold CV

- In K-Fold CV, the data are divided into K subsets

$$MSE_k = \frac{1}{n/K} \sum_{i \in S_k} (y_i - \hat{y}_i^{-S_k})^2$$

Cross-validation averages MSE

$$CV = \widehat{MSE} = \frac{1}{K} \sum_{k=1}^K MSE_k$$

- Each partition has a training dataset with $(K-1)*n/K$ observations and test dataset with n/K observations.
- Then the model is fit K times.
- LOOCV is a special case where the number of folds is n .

Model Selection

\widehat{g}^{-S_1} , (not S_1)
then apply
 \widehat{g}^{-S_1} to all
indices contained
in S_1 ?
all $i \in S_1$
obtain $\widehat{y}_i^{-S_1}$

Parametric Splines Model Selection

Training, validation, and test dataset

Issue: CV can still have issues w/ overfitting

Another approach is to divide the dataset into two datasets and use one dataset for estimation (training) and another for testing (test).

Popular splits: 70% (training) and 30% (test), or 80% and 20%.

This is often combined with cross-validation, where CV is applied to the training dataset to obtain the model, and then the model is applied to the unseen data to evaluate accuracy.

The "validation dataset" is the data left out when using cross-validation. In practice, we move across the different folds to tune hyperparameters. E.g., number of knots.

Popular in computer science with huge datasets.

Is there anything wrong with this approach? *requires a lot of data*

42/49 MS - Part 1

Parametric Splines Model Selection

K-Fold CV, cont. Bias-var tradeoff on MSE

- There is another bias-variance tradeoff ("meta" bias-variance tradeoff) regarding the optimal K
- K impacts the MSE of the sample estimate of the MSE,

$$E(MSE - \widehat{MSE})^2$$

where \widehat{MSE} is from K -fold CV.

- If we let $K = 2$, we will do a poorer job of fitting the model (using less data) \rightarrow upwardly biased \widehat{MSE} . *could underfit*
- At the other extreme, LOO has least bias because we are using $n - 1$ observations to estimate the models, but each term in the summand is highly correlated, leading to higher variance \rightarrow *higher variance of MSE*.
- General recommendations: 5 or 10-fold CV to balance this tradeoff.
- A discussion is James et al 2013 Chapter 5.

we will use for elasticnet k=5 or 10 fold for splines, we will use GCV or REML

43/49 MS - Part 1

Parametric Splines Model Selection

Likelihood-Based Approach

With basis expansion, we assume g_i can be expressed as a linear combination of some linear regressors,

$$\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2\mathbf{I}).$$

The Gaussian (Normal) likelihood is given by

$$f(\beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)\right\}$$

with log-likelihood function

$$\log f(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta).$$

Given our estimated coefficients $\hat{\beta}$, we hope to maximize

$$E[\log f(\hat{\beta}, \sigma^2)] = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} E[(\mathbf{X}\beta - \mathbf{X}\hat{\beta})'(\mathbf{X}\beta - \mathbf{X}\hat{\beta})].$$

44/49 MS - Part 1

Parametric Splines Model Selection

Mallow's C_p and Akaike Information Criterion

It can be shown that *for Gaussian w/ known σ^2*

$$E[\log f(\hat{\beta}, \sigma^2)] \approx \log f(\hat{\beta}, \sigma^2) + \frac{n}{2} - p.$$

45/49 MS - Part 1

It can be shown that for Gaussians w/ known σ^2

$$E[\log f(\hat{\beta}, \sigma^2)] \approx \log f(\hat{\beta}, \sigma^2) + \frac{n}{2} - p.$$

So one approach is to select a model that maximizes the above. If σ^2 is assumed known, this gives rise to the Mallows's C_p criterion:

$$C_p = \frac{1}{\sigma^2} (\mathbf{Y} - \hat{\beta}\mathbf{X})'(\mathbf{Y} - \hat{\beta}\mathbf{X}) - n + 2p$$

In practice, σ^2 is often replaced by its estimate from the largest model considered.

Akaike information criterion (AIC) generalizes this to any likelihood:

$$AIC = -2 \times \log f(\hat{\beta}, \hat{\sigma}^2) + 2p$$

which we wish to minimize also.

- can be viewed as an approximation of CV error (MSE)
for gaussian

43/49

Parametric Splines

Model Selection

Example: Mortality Trends

```
> ### AIC, CV, and GCV calculation examples:
> fit = lm(alldeaths~bs(date, df = 10), data = health)
> ##AIC
> AIC(fit)
[1] 15252.54
>
> ##GCV/CV
> H = hatvalues(fit)
> CV = mean((fit$fitted - health$alldeaths)^2 / (1-H)^2)
> CV
[1] 248.0374
>
> GCV = mean((fit$fitted - health$alldeaths)^2) / (1 - mean(H))^2
> GCV
[1] 248.1518
```

46/49

Parametric Splines

Model Selection

Example: Daily Death Counts in NYC

Optimal = 70 df. (the minimum)

AIC

CV

GCV

not monotonic

Death

2001 2002 2003 2004 2005 2006

47/49

