

Using splines in regression

Author: Nicholas G Reich, Jeff Goldsmith

*This material is part of the **statsTeachR** project*

*Made available under the Creative Commons Attribution-ShareAlike 3.0 Unported
License: http://creativecommons.org/licenses/by-sa/3.0/deed.en_US*

Today's Lecture

- Spline models
- Penalized spline regression

[More info: Harrel, *Regression Modeling Strategies*, Chapter 2, PDF handout]

Piecewise linear models

A piecewise linear model (also called a change point model or broken stick model) contains a few linear components

- Outcome is linear over full domain, but with a different slope at different points
- Points where relationship changes are referred to as “change points” or “knots”
- Often there's one (or a few) potential change points

Piecewise linear models

Suppose we want to estimate $E(y|x) = f(x)$ using a piecewise linear model.

- For one knot we can write this as

$$E(y|x) = \beta_0 + \beta_1 x + \beta_2 (x - \kappa)_+$$

where κ is the location of the change point and

$$(x - \kappa)_+ =$$

Interpretation of regression coefficients

$$E(y|x) = \beta_0 + \beta_1 x + \beta_2 (x - \kappa)_+$$

- $\beta_0 =$

- $\beta_1 =$

- $\beta_2 =$

- $\beta_1 + \beta_2 =$

Estimation

- Piecewise linear models are low-dimensional (no need for penalization)
- Parameters are estimated via OLS
- The design matrix is ...

Multiple knots

Suppose we want to estimate $E(y|x) = f(x)$ using a piecewise linear model.

- For multiple knots we can write this as

$$E(y|x) = \beta_0 + \beta_1 x + \sum_{k=1}^K \beta_{k+1} (x - \kappa_k)_+$$

where $\{\kappa_k\}_{k=1}^K$ are the locations of the change points

- Note that knot locations are defined before estimating regression coefficients
- Also, regression coefficients are interpreted conditional on the knots.

Example: lidar data

```
library(MASS)
library(SemiPar)

## Error:  there is no package called 'SemiPar'

data(lidar)

## Warning:  data set 'lidar' not found

y = lidar$logratio

## Error:  object 'lidar' not found

range = lidar$range

## Error:  object 'lidar' not found

qplot(range, y)

## Error:  object 'y' not found
```


Example: lidar data

```
knots <- c(550, 625)
mkSpline <- function(k, x) (x - k > 0) * (x - k)
X.des = cbind(1, range, sapply(knots, FUN = mkSpline, x = range))

## Error: non-numeric argument to binary operator

colnames(X.des) <- c("intercept", "range", "range1", "range2")

## Error: object 'X.des' not found

lm.lin = lm(y ~ X.des - 1)

## Error: object 'y' not found

plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)

## Error: object 'y' not found

points(range, lm.lin$fitted.values, type = "l", col = "red", lwd = 2)

## Error: object 'lm.lin' not found
```

Example: lidar data

```
summary(lm.lin)$coef
```

```
## Error: object 'lm.lin' not found
```

Piecewise quadratic and cubic models

Suppose we want to estimate $E(y|x) = f(x)$ using a piecewise quadratic model.

- For multiple knots we can write this as

$$E(y|x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \sum_{k=1}^K \beta_{k+2} (x - \kappa_k)_+^2$$

where $\{\kappa_k\}_{k=1}^K$ are the locations of the change points

- Similar extension for cubics
- Piecewise quadratic models are smooth and have continuous first derivatives
- Often, knots taken as quintiles of the data.

Advantages of piecewise models

Piecewise (linear, quadratic, etc) models have several advantages

- Easy construction of basis functions
- Flexible, and don't rely on determining an appropriate form for $f(x)$ using standard functions
- Allow for significance testing on change point slopes
- Fairly direct interpretations

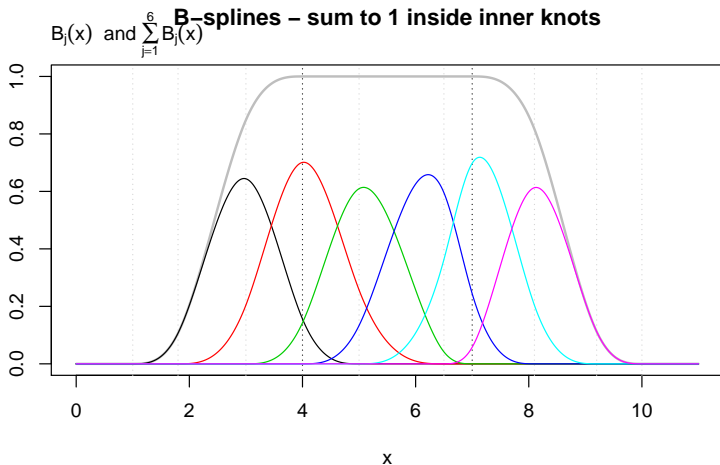
B-splines and natural splines

Characteristics

- Both B-splines and natural splines similarly define a basis over the domain of x
- Can be constrained to have seasonal patterns
- They are made up of piecewise polynomials of a given degree, and have defined derivatives similarly to the piecewise defined functions
- Big advantage over linear splines: parameter estimation is often fairly robust to your choice of knots
- Big disadvantage over linear splines: harder to interpret specific coefficients

B-splines basis functions

$$E(y|x) = \beta_0 + \sum_{j=1}^6 \beta_j B_j(x)$$



Example: lidar data

```
require(splines)
lm.bs3 = lm(y ~ bs(range, df = 3))

## Error: object 'y' not found

plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)

## Error: object 'y' not found

points(range, lm.bs3$fitted.values, type = "l", col = "red", lwd = 2)

## Error: object 'lm.bs3' not found
```

Example: lidar data

```
lm.bs5 = lm(y ~ bs(range, df = 5))  
  
## Error: object 'y' not found  
  
plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)  
  
## Error: object 'y' not found  
  
points(range, lm.bs5$fitted.values, type = "l", col = "red", lwd = 2)  
  
## Error: object 'lm.bs5' not found
```