

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 —

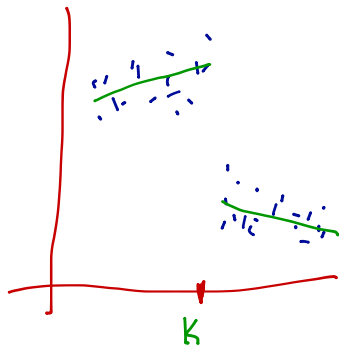
More info:

- Harrel, *Regression Modeling Strategies*, Chapter 2, PDF handout
- *ISL* Chapter 7

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



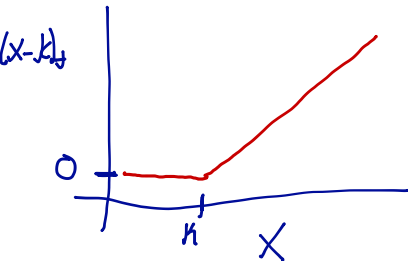
Piecewise linear models

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



$$(x - \kappa)_+ = \begin{cases} x - \kappa, & x > \kappa \\ 0, & \text{ow.} \end{cases}$$
$$= \max(x - \kappa, 0)$$

Interpretation of regression coefficients

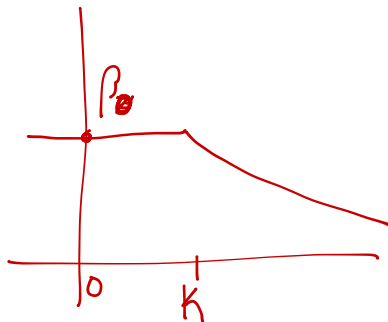
$$E(y|x) = \beta_0 + \beta_1 x + \beta_2 (x - \kappa)_+$$

■ $\beta_0 = \mathbb{E}[y|x = 0]$ (assuming $\kappa > 0$)

■ $\beta_1 =$

■ $\beta_2 =$

■ $\beta_1 + \beta_2 =$

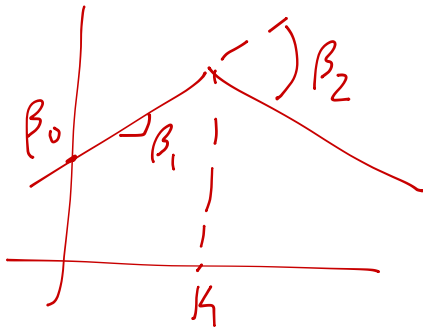


Interpretation of regression coefficients

$x < \kappa$

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

- $\beta_0 = \mathbb{E}[y|x = 0]$ (assuming $\kappa < 0$)
- $\beta_1 =$ Expected change in y for a 1-unit increase in x , when $x < \kappa$
- $\beta_2 =$
- $\beta_1 + \beta_2 =$



Interpretation of regression coefficients

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

- $\beta_0 = \mathbb{E}[y|x = 0]$ (assuming $\kappa < 0$)
- $\beta_1 =$ Expected change in y for a 1-unit increase in x , when $x < \kappa$
- $\beta_2 =$ Change in slope between $x < \kappa$ and $x > \kappa$
- $\beta_1 + \beta_2 =$

Interpretation of regression coefficients

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

- $\beta_0 = \mathbb{E}[y|x = 0]$ (assuming $\kappa < 0$)
- $\beta_1 =$ Expected change in y for a 1-unit increase in x , when $x < \kappa$
- $\beta_2 =$ Change in slope between $x < \kappa$ and $x > \kappa$
- $\beta_1 + \beta_2 =$ Expected change in y for a 1-unit increase in x , when $x \geq \kappa$

Estimation

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

$\beta_0 \quad \beta_1 \quad \beta_2 \quad \leftarrow (X-K)_k$

$$X = \begin{bmatrix} 1 & X_1 & 0 \\ 1 & X_2 & 0 \\ 1 & X_3 & 0 \\ 1 & X_4 & 0 \\ \vdots & \vdots & X-K \\ \vdots & \vdots & X-K \\ 1 & \vdots & \end{bmatrix}$$

low
high

Sort by X

Multiple knots

$$f(x|K) = (x-K)_+ \\ = \max(x, 0) \cdot x - K$$

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

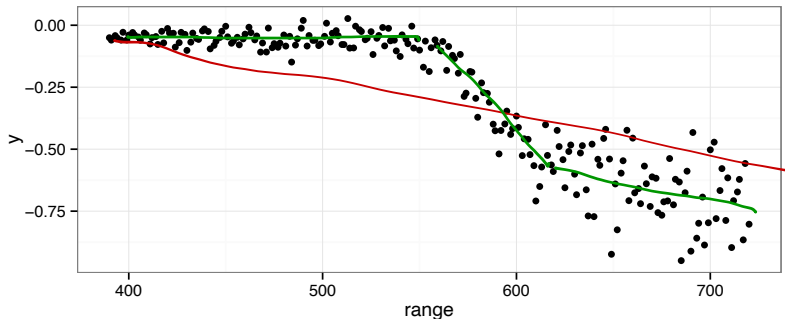
- Note that knot locations are defined before estimating regression coefficients
- Also, regression coefficients are interpreted conditional on the knots.

Example: lidar data

```
library(MASS)
library(SemiPar)
```

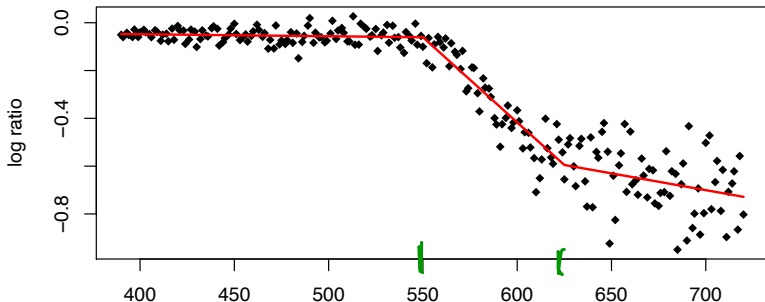
```
## Warning: package 'SemiPar' was built under R version 3.1.2
```

```
data(lidar)
y = lidar$logratio
range = lidar$range
qplot(range, y)
```



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))
colnames(X.des) <- c("intercept", "range", "range1", "range2")
lm.lin = lm(y ~ X.des - 1)
plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)
points(range, lm.lin$fitted.values, type = 'l', col = "red", lwd = 2)
```



Example: lidar data

lincom()

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

##		Estimate	Std. Error	t value	Pr(> t)
##	X.desintercept	-1.444288e-02	0.0687353855	-0.2101230	8.337689e-01
##	X.desrange	β_1 -8.407376e-05	0.0001426647	-0.5893102	5.562663e-01
##	X.desrange1	β_2 -7.042794e-03	0.0003834218	-18.3682689	4.379404e-46
##	X.desrange2	β_3 5.723186e-03	0.0005153479	11.1054811	5.554824e-23

$\beta_1 = \text{slope of 2}^{\text{nd}} \text{ line}$
 $\beta_1 + \beta_2 = \text{slope of 2}^{\text{nd}}$
 $\beta_1 + \beta_2 + \beta_3 = \text{slope of 3}^{\text{rd}}$

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

Pros and cons 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 *esp. for quadratic*
- Allow for significance testing on change point slopes
- Fairly direct interpretations

Disadvantages

- knot specification is often arbitrary

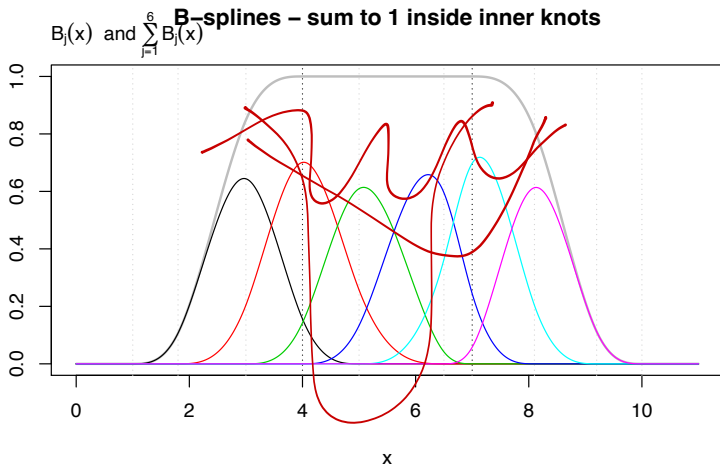
B-splines and natural splines

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

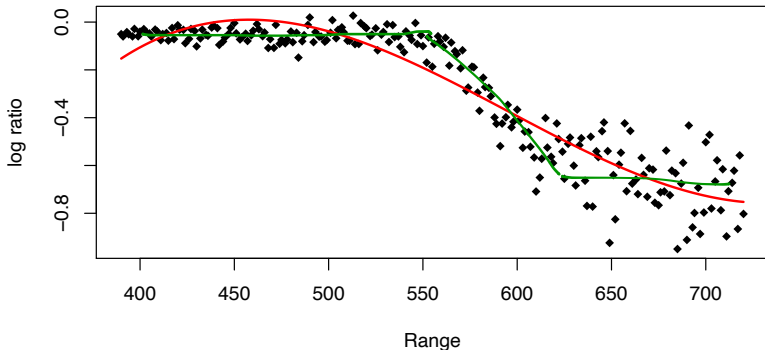
$$E(y|x) = \beta_0 + \sum_{j=1}^6 \beta_j B_j(x)$$



Example: lidar data

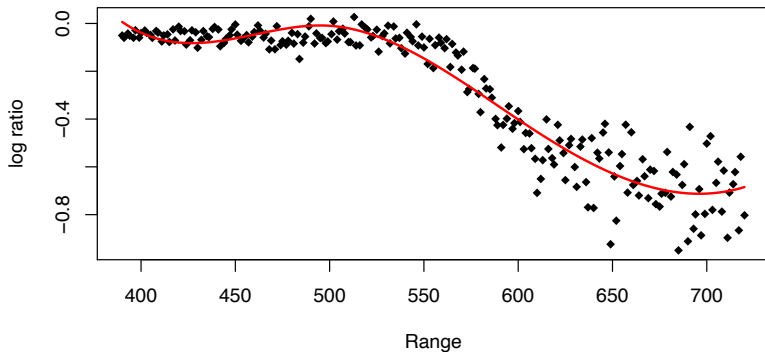
```
require(splines)
lm.bs3 = lm(y ~ bs(range, df=3))
plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)
points(range, lm.bs3$fitted.values, type = 'l', col = "red", lwd = 2)
```

→ knots



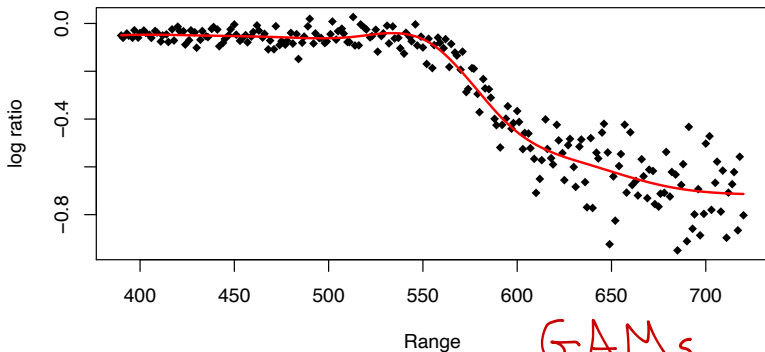
Example: lidar data

```
lm.bs5 = lm(y ~ bs(range, df=5))  
plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)  
points(range, lm.bs5$fitted.values, type = 'l', col = "red", lwd = 2)
```



Example: lidar data

```
lm.bs5 = lm(y ~ bs(range, df=10))  
plot(range, y, xlab = "Range", ylab = "log ratio", pch = 18)  
points(range, lm.bs5$fitted.values, type = 'l', col = "red", lwd = 2)
```



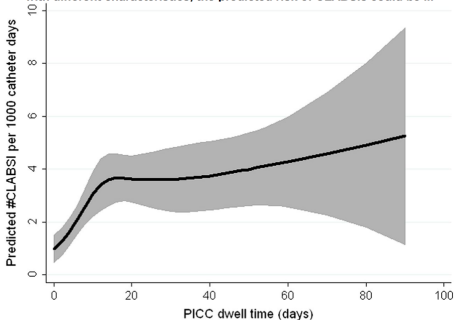
GAMs
maew

Take-home points for spline approaches (1)

Spines can flexibly model non-linear relationships

- Can improve model fit because of relaxed linearity assumptions.
- Caveat: spline models require careful graphical interpretation, slopes may not be easily available/interpretable

Predicted CLABSI rate (solid line) with 95% CIs (shading) over catheter dwell time for a given hospital, age, birth weight, CLABSI from previous PICC, and concurrent PICC. For a neonate with different characteristics, the predicted risk of CLABSIs could be ...



Take-home points for spline approaches (2)

Do you want control over your knots?

- Your application may have explicit “change-points” (i.e. interrupted time-series)
- In most cases, you do not want your spline model to be sensitive to user input (i.e. knot placement)
- “Penalized splines” can reduce this sensitivity at the cost of more complex model and estimation (More in *ISL* Chapter 7, Biostat Methods 3, anything about Generalized Additive Models (e.g. `mgcv` package and `gam()` function), one of your projects?).