

Lectures 1, 2, and 4

Splines, Smoothers, and Generalized Additive Models

Data Science 2
CS 109b, Stat 121b, AC 209b, E-109b

Mark Glickman Pavlos Protopapas Chris Tanner

Reading: James et al., chapter 7.

Outline

- Polynomial regression and basis functions
- Regression splines
- Smoothers
- Additive and Generalized Additive Models

From linear to non-linear effects

You have seen in CS 109a models in which the contribution of a predictor is included linearly.

For a quantitative response Y and quantitative predictor x , we can assume

$$Y = \beta_0 + \beta_1 x + \varepsilon$$

with $\varepsilon \sim N(0, \sigma^2)$.

We have a number of options for including non-linear effects of x if we believe the relationship is not linear.

Some methods you have already seen in CS 109a (e.g., random forests).

Review – polynomial models in one variable

(there will be an ultimate point to this review – be patient!)

Create new variables x^2, x^3, \dots, x^d for some specified d . Then assume

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d + \varepsilon.$$

Modeling a polynomial function of a predictor is useful

- when the researcher knows based on scientific theory that the true mean function is curvilinear with respect to the predictor.

- as an approximation to a possibly complex non-linear mean function of the predictor.

Comments on polynomial functions of predictors

- The model

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_d x^d + \varepsilon$$

is actually a linear regression model. This is because the model is linear in the unknown parameters β_j .

- Can use the same polynomial approach for logistic regression:

$$\text{logit } \Pr(Y = 1) = \log \left(\frac{\Pr(Y = 1)}{\Pr(Y = 0)} \right) = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_d x^d.$$

- Can use standard model fitting methods (least-squares, maximum likelihood) for polynomial models because of the linearity of the β_j .

Why polynomial models make sense: Taylor's theorem

Suppose $f(x)$ is the true mean function of x , so that

$$E(Y|x) = f(x).$$

Taylor's theorem says (in essence) that we can approximate $f(x)$ to any degree of accuracy we like with a polynomial approximation, with better accuracy for higher order polynomials.

$$\begin{aligned} f(x) &= f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \cdots + \frac{f^{(d)}(0)}{d!}x^d + \cdots \\ &= \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_d x^d + \cdots \end{aligned}$$

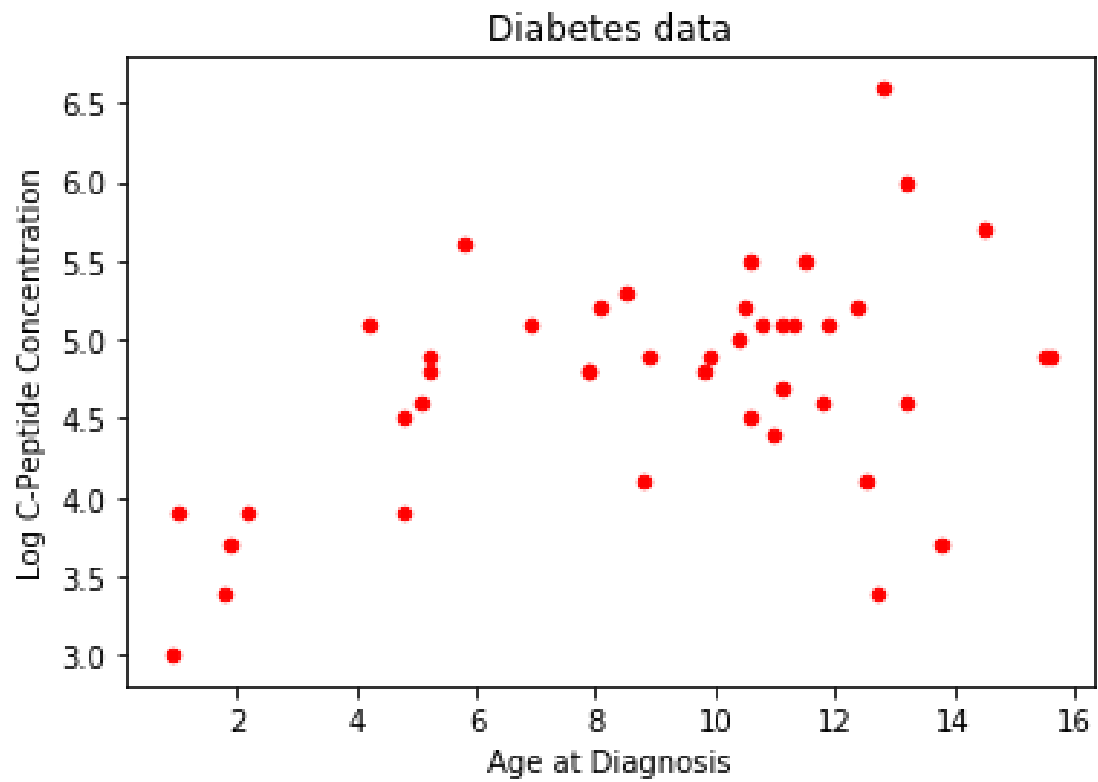
Simple motivating example: Diabetes in children

A study from 1987 investigated factors affecting diabetes in children as measured through serum C-peptide levels.

- Response/Label: $Y = \text{Log of serum C-peptide level (pmol/ml) in a child}$
- Predictor/Feature: $x = \text{Age the child was diagnosed with diabetes}$

Want to examine the relationship between these two variables on the 43 children in the study.

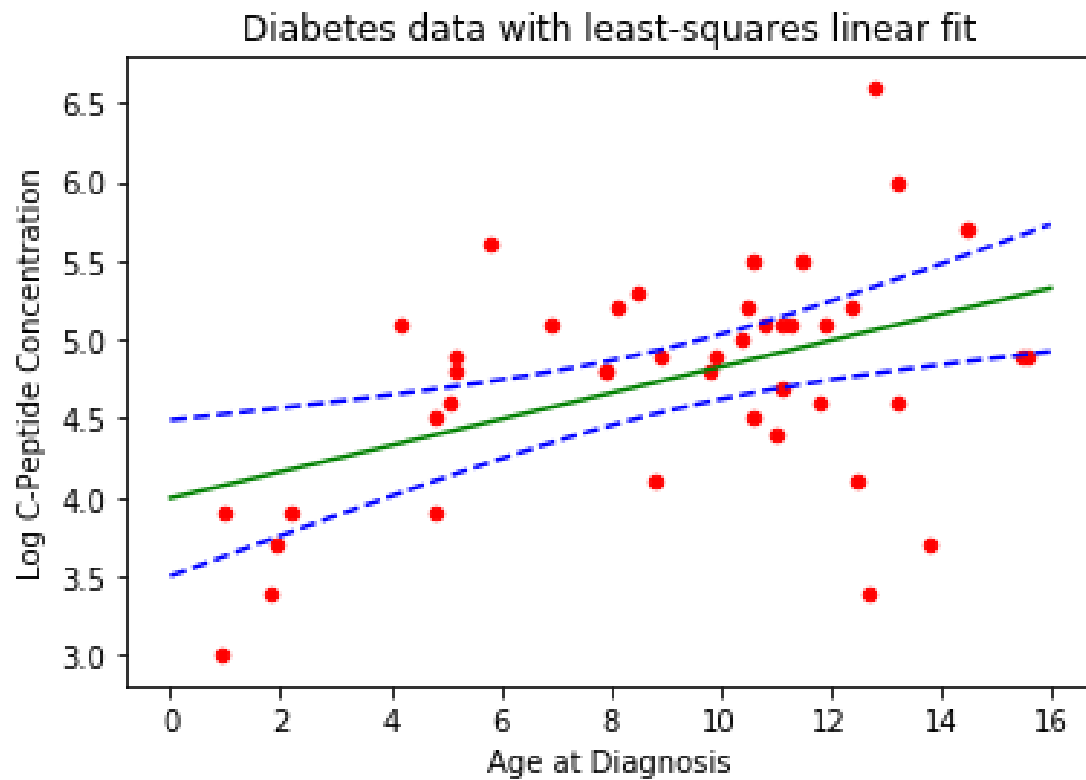
Diabetes in children



[First try: Least squares regression line](#)

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

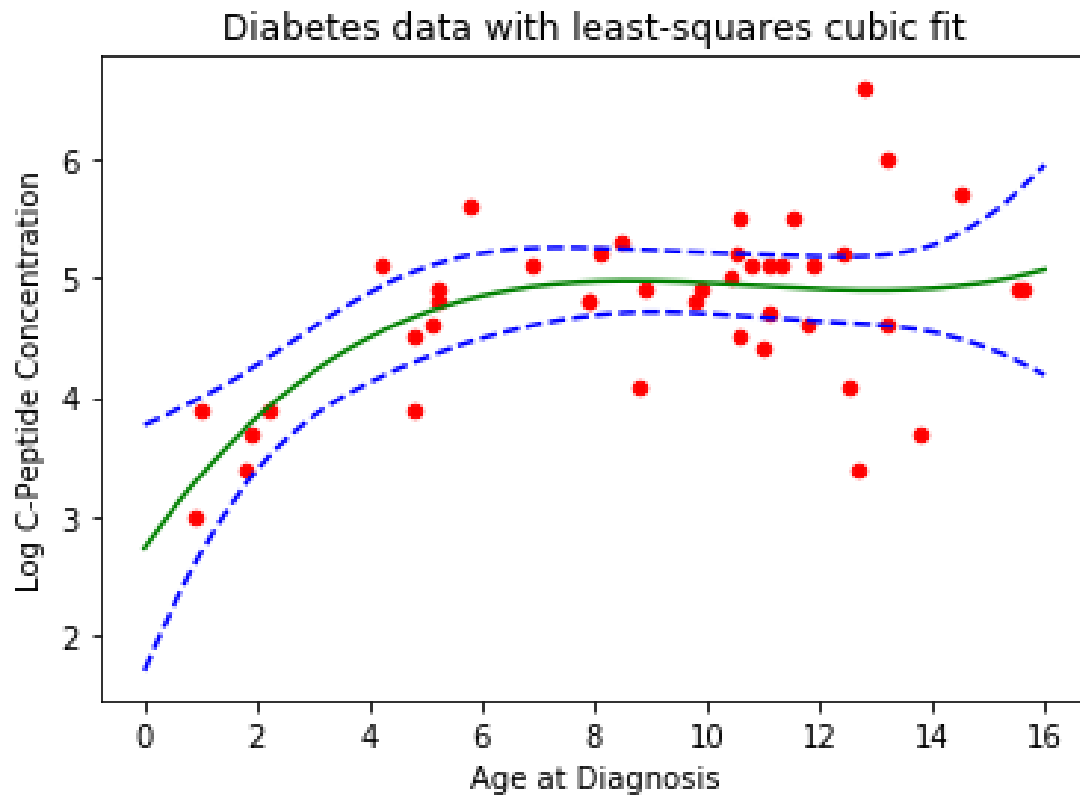
[Diabetes in children](#)



[Second try: Least squares regression cubic polynomial](#)

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

[Diabetes in children](#)



Can fit logistic regression models as well.

Define

$$Y_i^* = \begin{cases} 1 & \text{if log C-Peptide concentration} > 4.0 \\ 0 & \text{if log C-Peptide concentration} \leq 4.0 \end{cases}$$

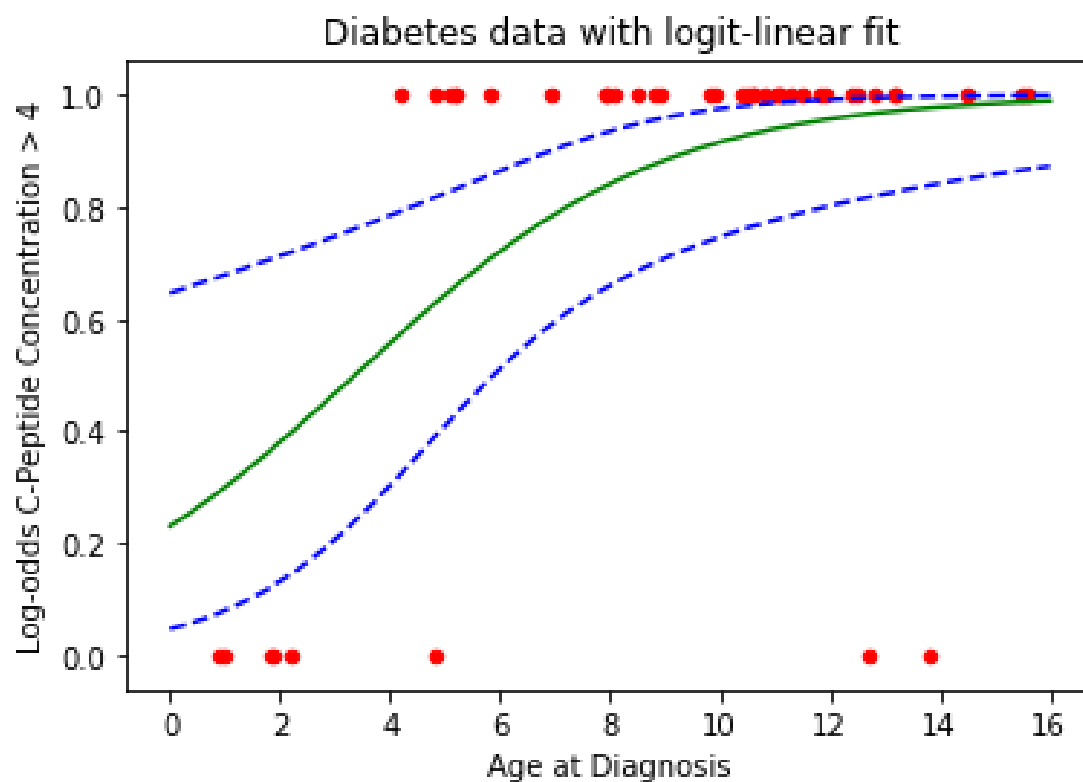
Linear age:

$$\text{logit } \Pr(Y_i^* = 1) = \beta_0 + \beta_1 x_i.$$

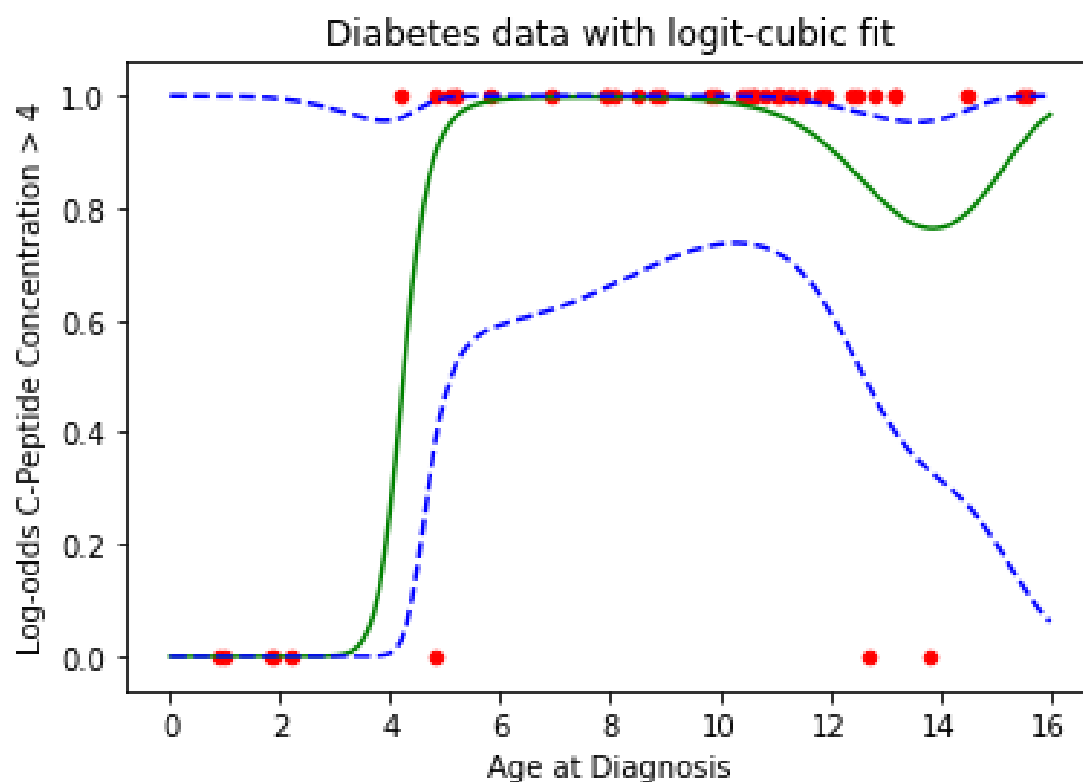
Cubic age:

$$\text{logit } \Pr(Y_i^* = 1) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3.$$

Diabetes in children



[Diabetes in children](#)



Comments

- Can include polynomial terms for multiple variables simultaneously in multiple least-squares or logistic regression.
- Choice of d (order of polynomial) - either by scientific considerations, or by cross-validation.
- Actually rare to model unknown curvature using polynomial functions
 - Usually low-order polynomials provide a poor fit
 - Larger order polynomials do not capture local behavior well unless d is very large.
 - Large d can be difficult to interpret.

Deeper look into nonlinear regression

Goal Model $E(Y|x)$ as a function of a predictor, x , based on a flexible set of choices of functions.

The true model may be

$$E(Y|x) = f(x)$$

where $f(x)$ is a highly nonlinear function of predictor variable x .

We may be willing to consider as best approximations:

$$E(Y|x) = \beta_0 + \beta_1 x$$

or

$$E(Y|x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_d x^d.$$

In each case, we are considering linear combinations of component functions.

For linear regression, we have two functions of x :

$$\begin{aligned} f_1(x) &= 1 \\ f_2(x) &= x \end{aligned}$$

such that an approximation to $f(x)$ must be a linear combination of these two functions:

$$E(Y|x) = c_1 f_1(x) + c_2 f_2(x).$$

Similarly, for polynomial regression, we have

$$\begin{aligned} f_1(x) &= 1 \\ f_2(x) &= x \\ &\vdots \\ f_{d+1}(x) &= x^d \end{aligned}$$

so that $f(x)$ is approximated by the linear combination

$$E(Y|x) = c_1 f_1(x) + \cdots + c_{d+1} f_{d+1}(x).$$

Basis, and basis functions The set of component functions that are part of the approximating linear combination are called a basis, and the functions themselves are the corresponding basis functions.

The span of a basis is the set of all possible functions that can be formed from the linear combination of basis functions.

Examples

The set $S = \{1, x\}$ is a basis for all linear functions of x .

The set $S = \{1, x, x^2, \dots, x^d\}$ is a basis for all d -th order polynomials of x .

The point of the discussion of polynomial regression

- Polynomial regression is arguably the simplest example of creating a basis to approximate a non-linear function.
- Takes advantage of a linear combination representation

$$E(Y|x) = c_1 f_1(x) + \cdots + c_{d+1} f_{d+1}(x).$$

where each of the $f_j(x)$ are individual basis functions, namely $f_j(x) = x^{j-1}$.

Slightly more interesting examples

The (infinite) set $S = \{1, x, x^2, x^3, \dots\}$ is a basis for all differentiable functions of x . (by Taylor's theorem)

The (infinite) set

$$S = \{1, \sin(2\pi x), \sin(4\pi x), \sin(6\pi x), \dots, \cos(2\pi x), \cos(4\pi x), \cos(6\pi x), \dots\}$$

is a basis for all square-integrable functions of x , for $0 < x < 1$.

Orthogonal polynomial basis

We have already seen the use of a polynomial basis $\{1, x, x^2, \dots, x^d\}$.

One issue that can arise in a linear models setting is that terms like x and x^2 can be highly correlated for observed predictor data.

Example Suppose we have a data set of $n = 5$ values, and for predictor x we observe

$$x = (0.98, 0.99, 1.0, 1.01, 1.02).$$

Then the values of x^2 are

$$x^2 = (0.9604, 0.9801, 1.0, 1.0201, 1.0404).$$

The correlation between the vectors x and x^2 is 0.9999825.

[Problem](#) Having highly correlated predictors can

- result in numerical instability in fitting the model
- produce meaningless (highly inflated) coefficient estimates

[A solution](#) Instead of using $\{1, x, x^2\}$ as the basis, use a basis of polynomials that are uncorrelated on the data.

For the 5-observation example, let

$$S = \{1, (x - 1)/0.03152, -0.5345 + 2672.61(x - 1)^2\}$$

This basis is an example of an orthogonal polynomial basis.

Orthogonal polynomial basis

- This basis is comprised of a constant, a linear function of x , and a quadratic function of x . It forms a basis for all quadratic functions just like $\{1, x, x^2\}$.
- The linear and quadratic basis functions for the five observations evaluates to

$$(x - 1)/0.03152 = (-0.632, -0.316, 0.000, 0.316, 0.632)$$

and

$$-0.5345 + 2672.61(x - 1)^2 = (0.535, -0.267, -0.535, -0.267, 0.535)$$

The correlation between these two vectors is 0. Also, the correlation between each vector and a column of 1s is also 0.

The values are normalized so that the sum of the squared elements is 1.

Back to polynomial regression

Polynomial regression is actually rarely used in practice unless a scientific theory dictates the use of polynomials.

- Low-order polynomials are an inferior solution to other existing ways to acknowledge non-linearity.
- Increasing the order of the polynomial does not usually help because of odd behavior near the extremes of the data.

A more flexible approach is to use piecewise polynomials.

In particular, use connected piecewise polynomials, also known as splines.

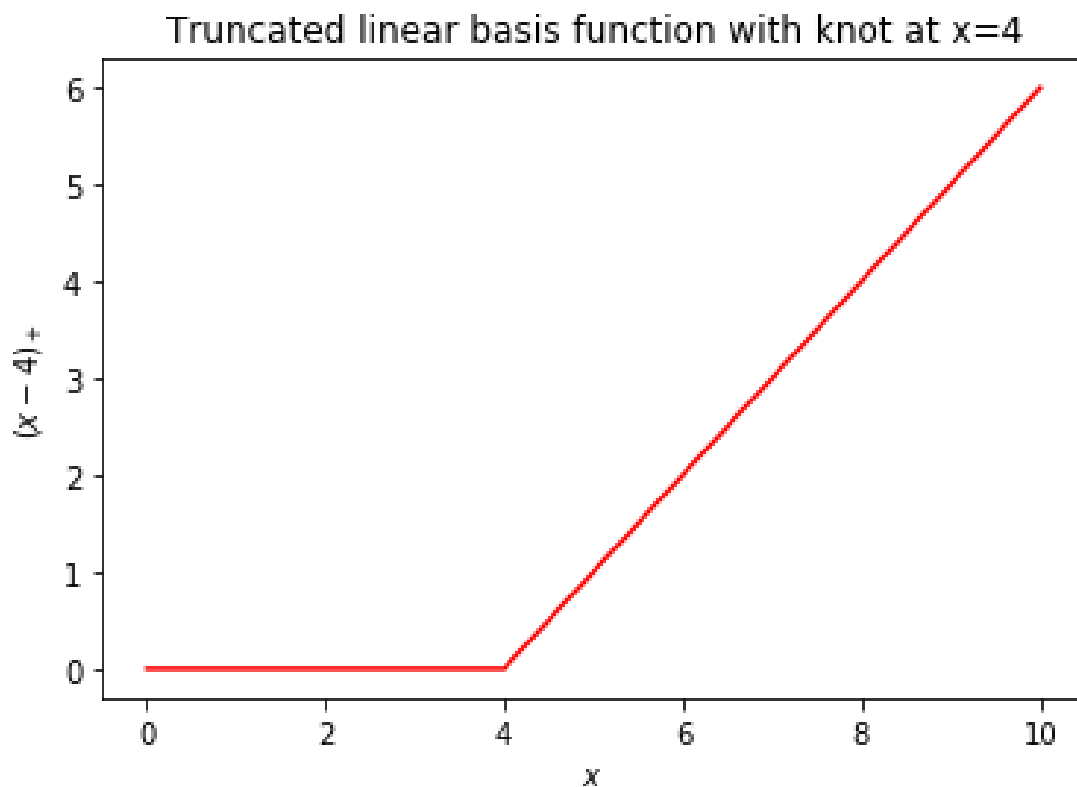
Example: Piecewise linear spline

Define the function

$$(x - \xi)_+ = \begin{cases} x - \xi & \text{if } x > \xi \\ 0 & \text{otherwise.} \end{cases}$$

This function is flat to the left of ξ , and linear (with slope 1) to the right. The value ξ is called a “knot” of the function.

The function $(x - \xi)_+$, a *truncated linear function*, will be a building block for constructing splines.



Connection to ReLUs

ReLU is short for “**R**ectified **L**inear **U**nit.”

A ReLU is a truncated linear function evaluated at 0, that is,

$$x_+ = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise.} \end{cases}$$

These functions are an essential building block for deep neural nets, as you started seeing last semester.

Strategy for a piecewise linear spline construction

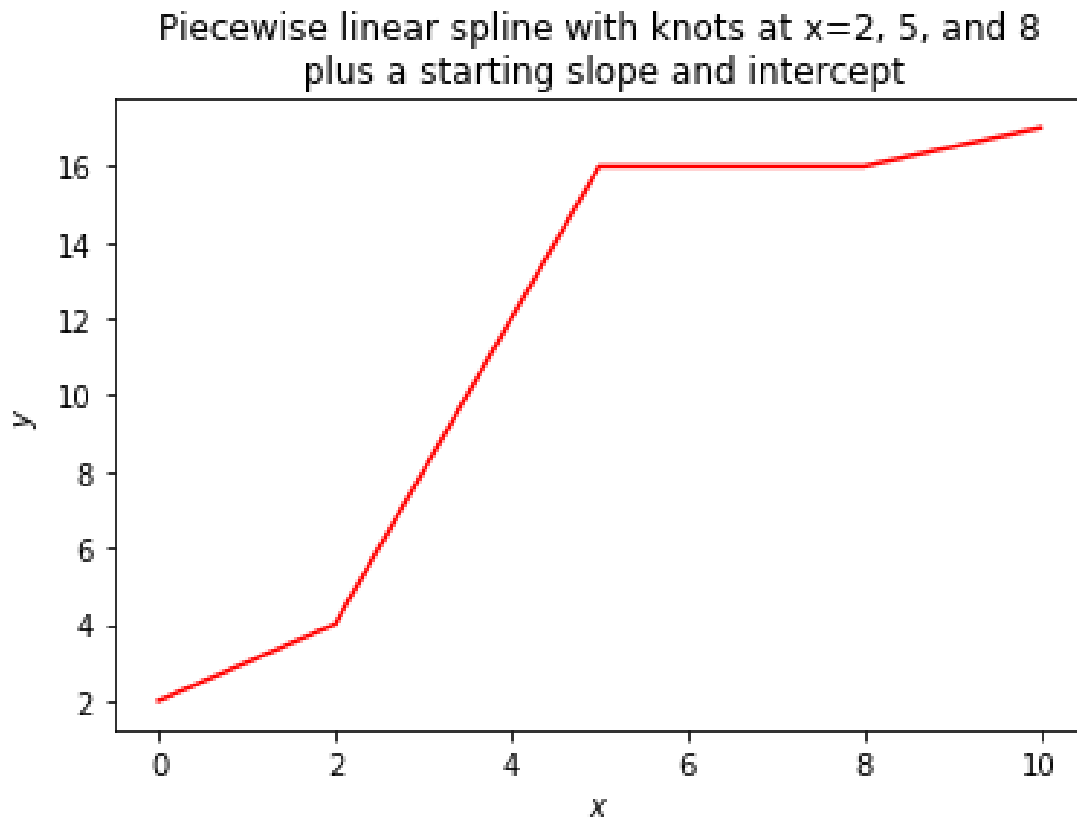
Choose knots $\xi_1 < \xi_2, \dots < \xi_K$

Now let

$$E(Y|x) = (\alpha_0 + \alpha_1 x) + \{\beta_1(x - \xi_1)_+ + \beta_2(x - \xi_2)_+ + \dots + \beta_K(x - \xi_K)_+\}.$$

Can think of this construction as

- Start with a linear function $(\alpha_0 + \alpha_1 x)$ to the left of ξ_1 .
- At each knot ξ_k , change the slope of the line by an additive amount β_k .



Spline basis functions

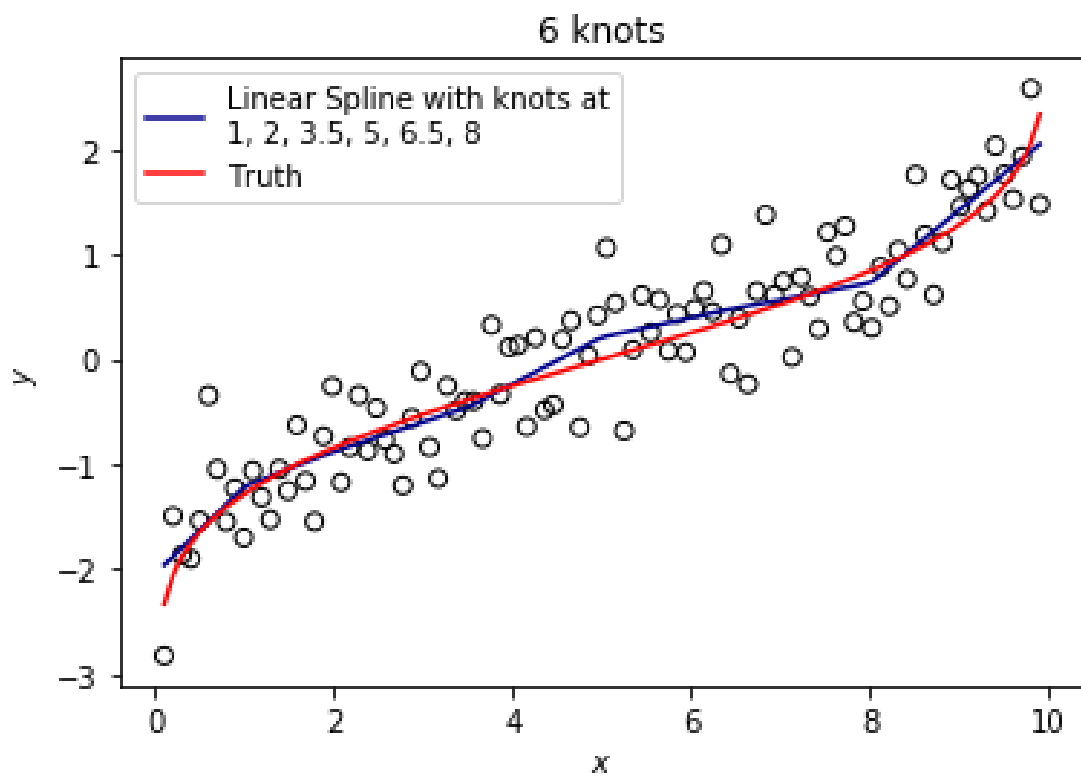
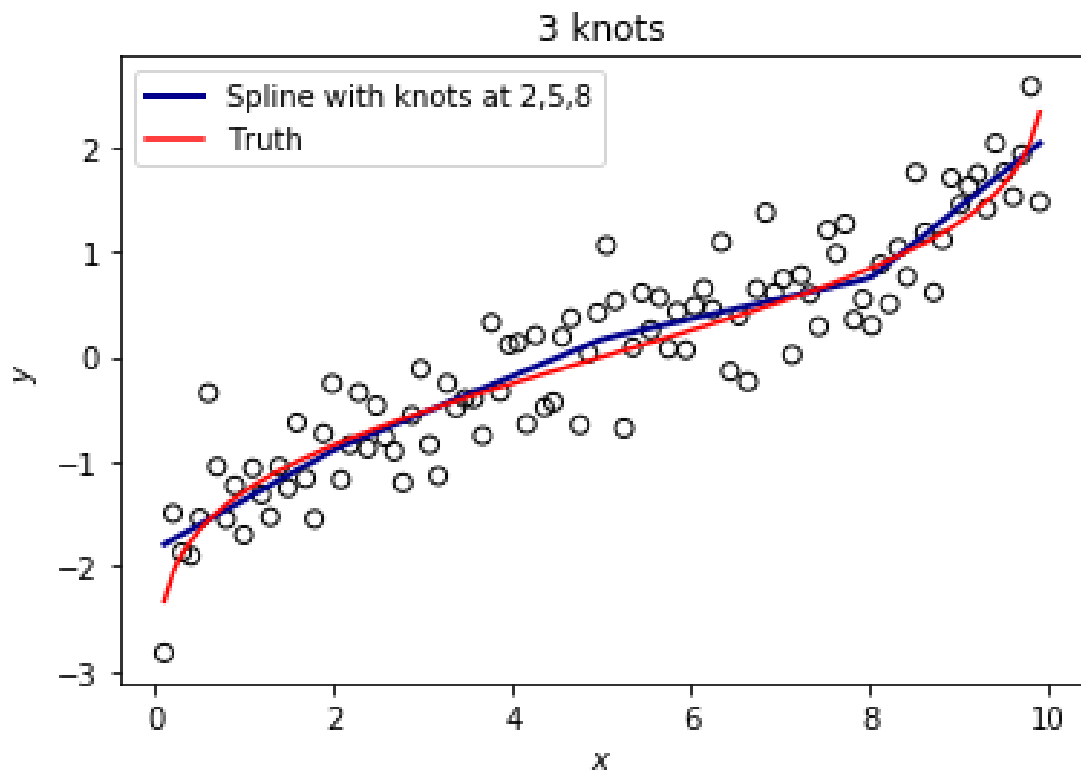
Given knots ξ_1, \dots, ξ_K , any linear spline can be composed from the following basis functions:

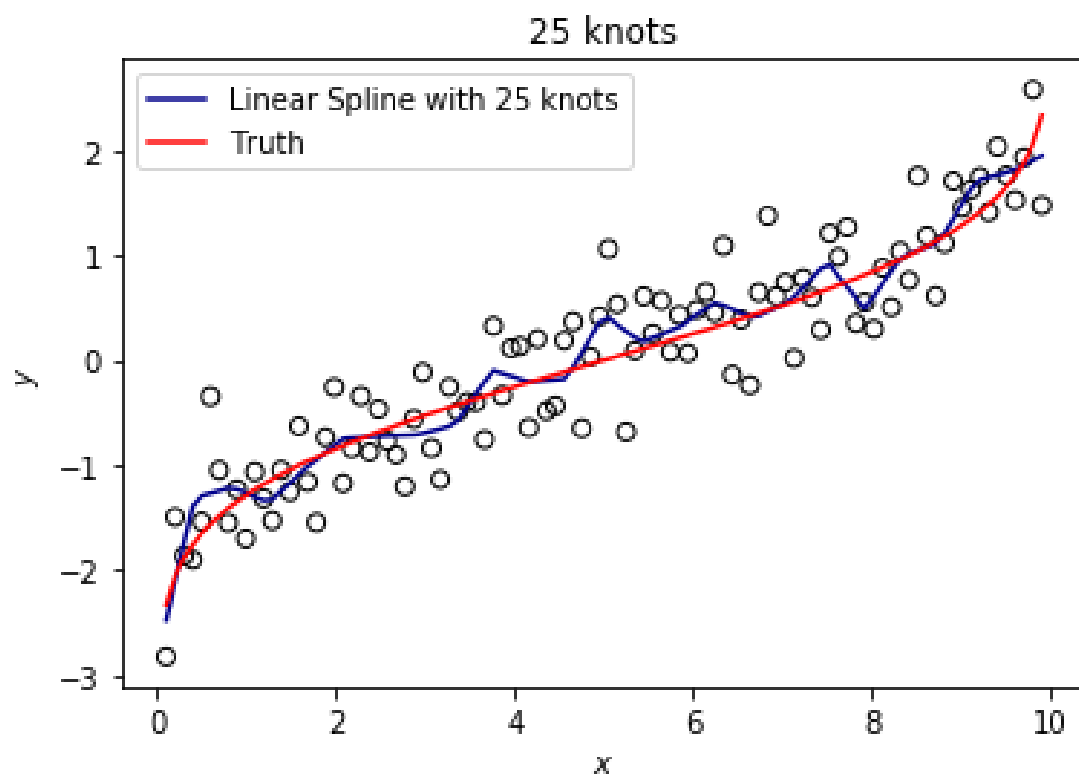
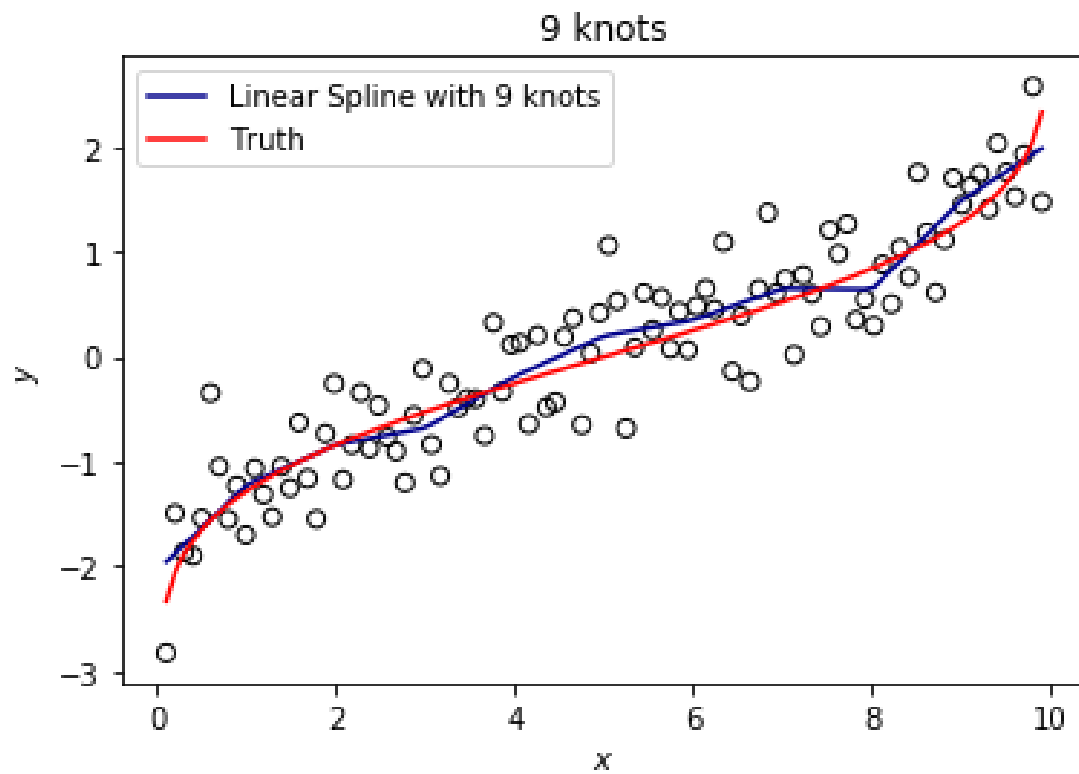
$$S = \{1, x, (x - \xi_1)_+, (x - \xi_2)_+, \dots, (x - \xi_K)_+\}$$

That is, S forms a basis for linear splines.

Increasing the number of knots

- increases the number of basis function components, and
- increases the scope of functions representable by a linear spline





Cubic splines

Linear splines are okay, but

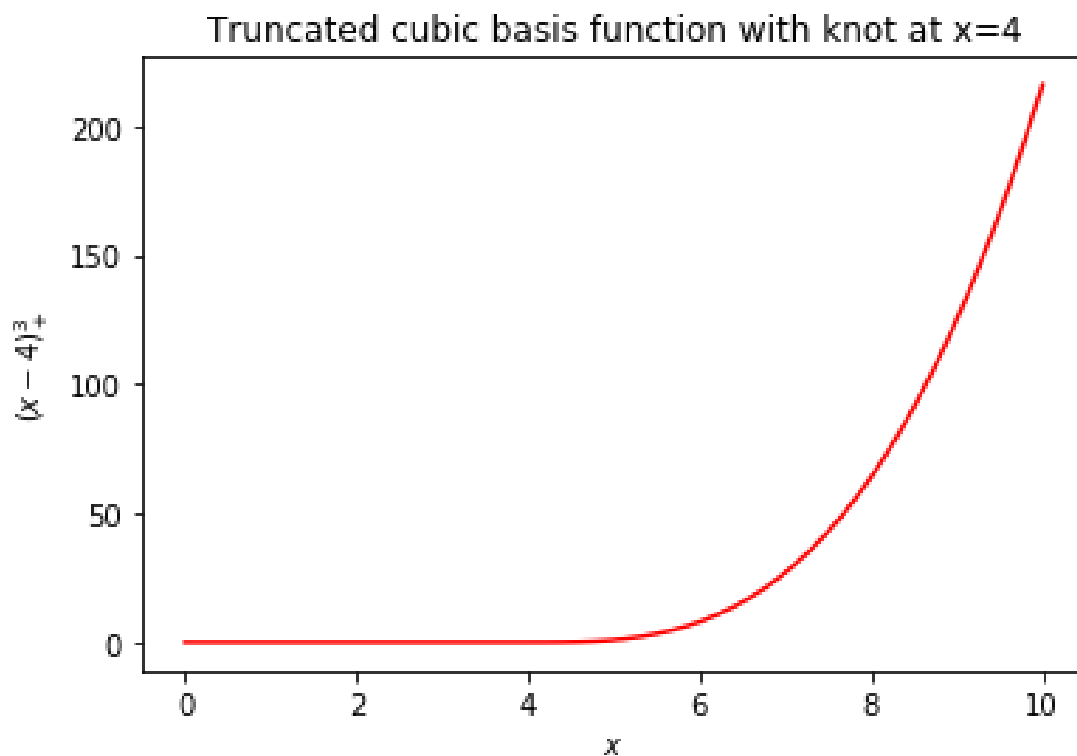
- they are not smooth
- they do not perform well detecting highly curved relationships (unless using a lot of knots)

Much more common in practice to use cubic splines.

The basis function that will help to develop cubic splines is

$$(x - \xi)_+^3 = \begin{cases} (x - \xi)^3 & \text{if } x > \xi \\ 0 & \text{otherwise.} \end{cases}$$

This is called a truncated cubic function.



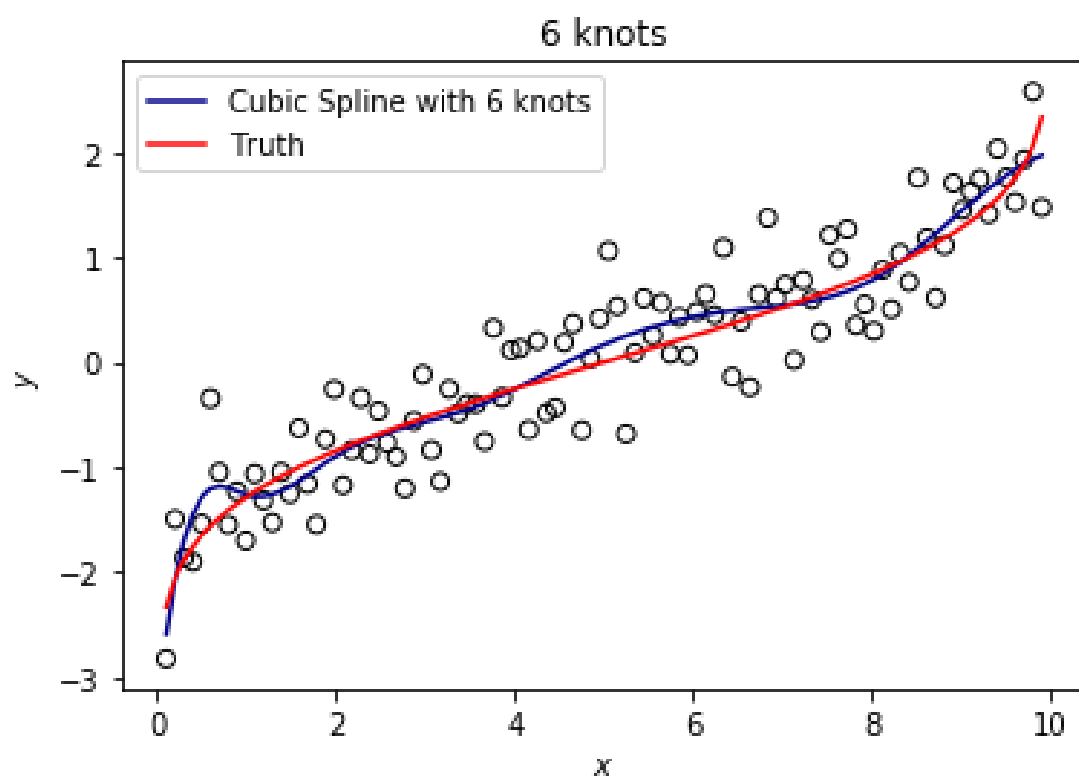
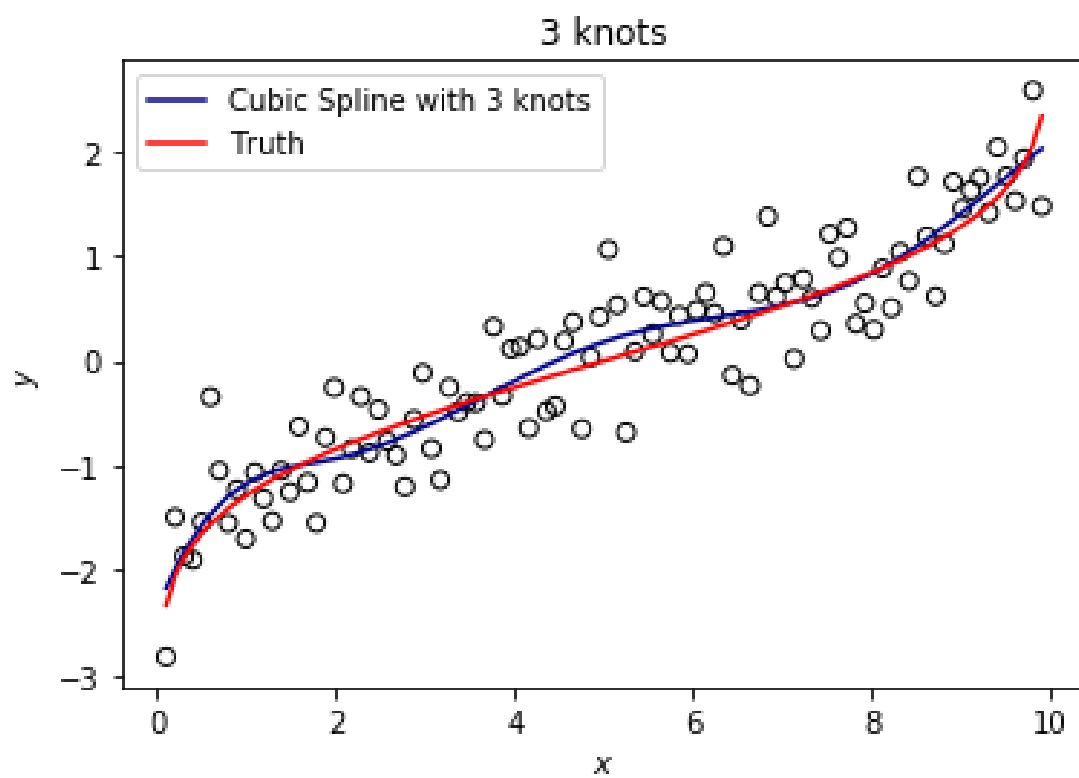
A cubic spline with K knots ξ_1, \dots, ξ_K to express the mean function of Y is given by

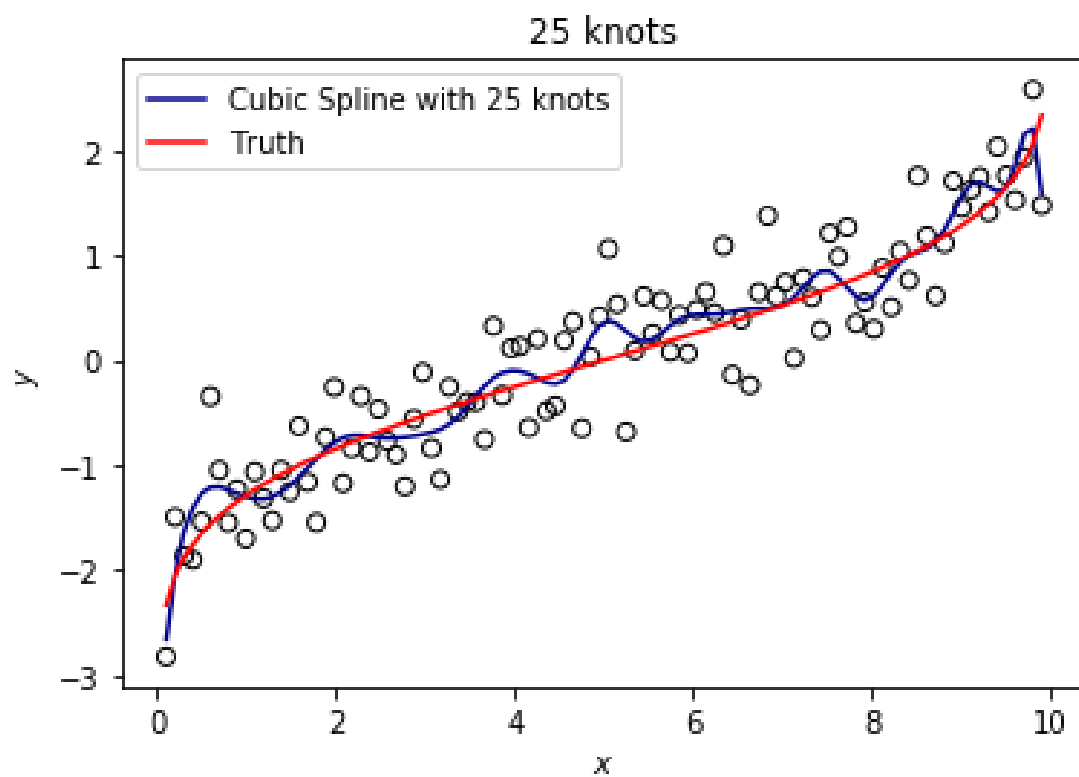
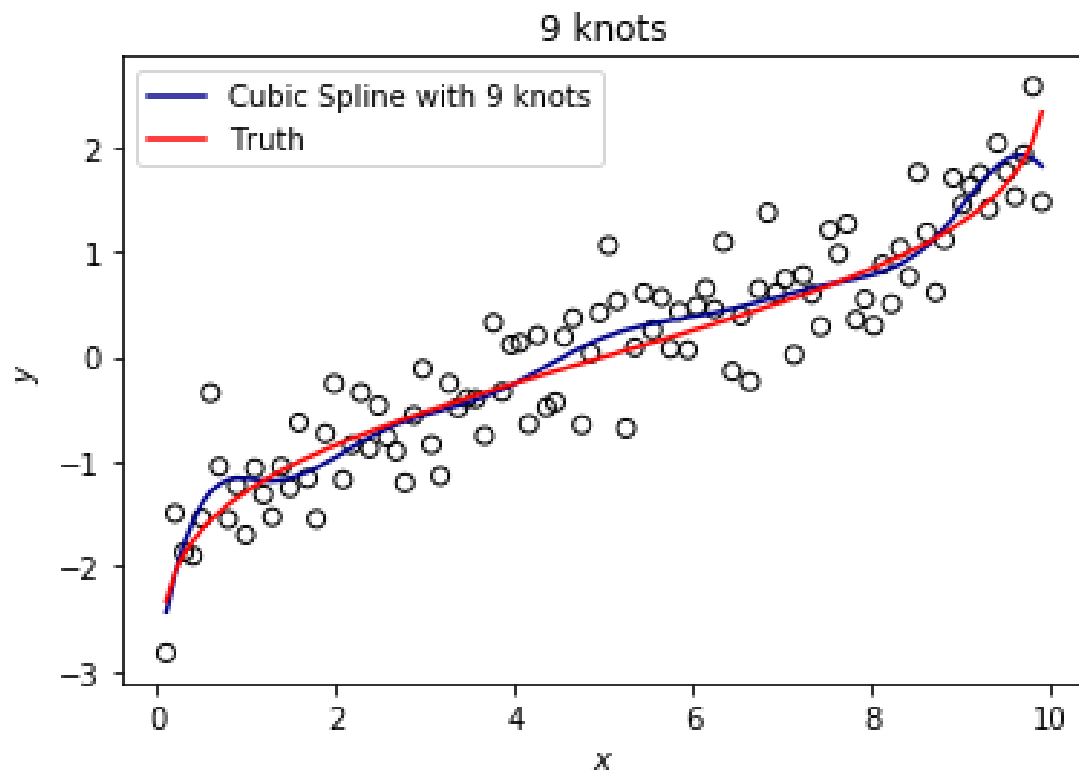
$$f_{cs}(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \alpha_3 x^3 + \{\beta_1(x - \xi_1)_+^3 + \beta_2(x - \xi_2)_+^3 + \dots + \beta_K(x - \xi_K)_+^3\}$$

Cubic spline basis: $S = \{1, x, x^2, x^3, (x - \xi_1)_+^3, \dots, (x - \xi_K)_+^3\}$

For quantitative outcomes, can perform least-squares regression on

$$y_i = f_{cs}(x_i) + \varepsilon_i$$





Important features of cubic spline

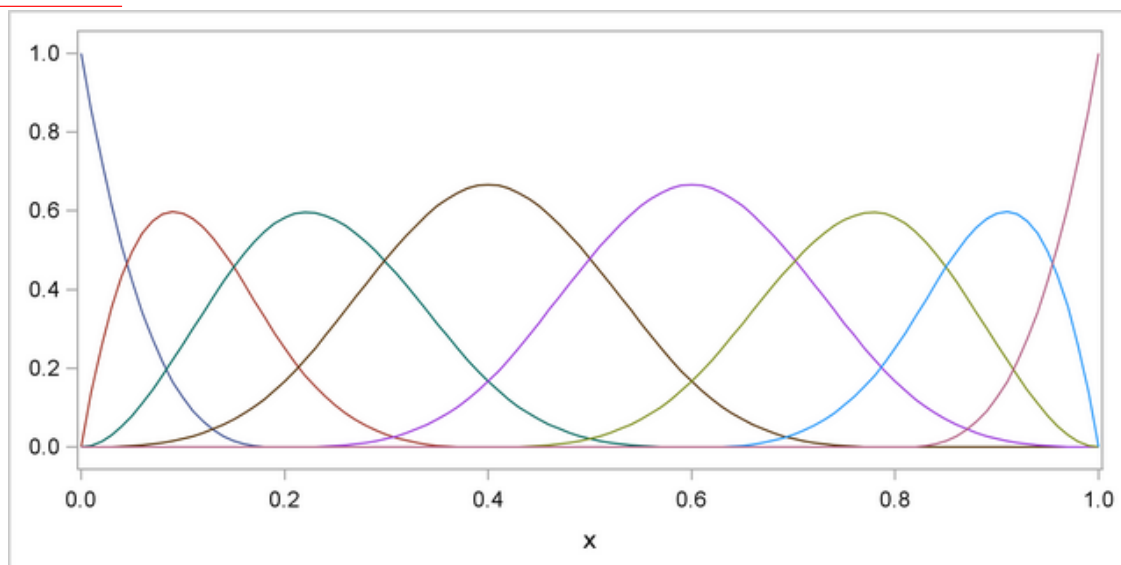
- In between knots ξ_k and ξ_{k+1} the spline as a function of x is a cubic polynomial. This is because the spline value is a sum of cubic polynomials (the “base” polynomial plus the truncated cubic terms to the left of ξ_k).
- At each knot, the spline is continuous, has a continuous first derivative, and has a continuous second derivative. This means that the spline is smooth (to the eye).
- A cubic spline, without any extra constraints, has $K+4$ parameters that need to be estimated.

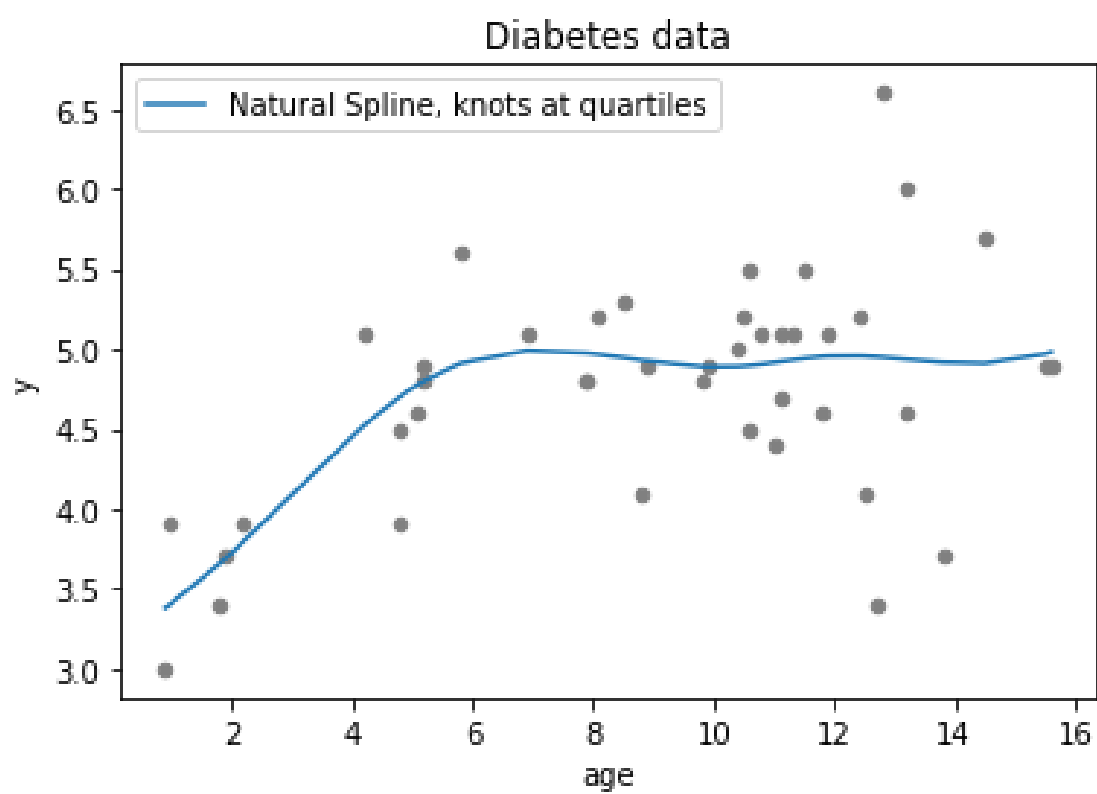
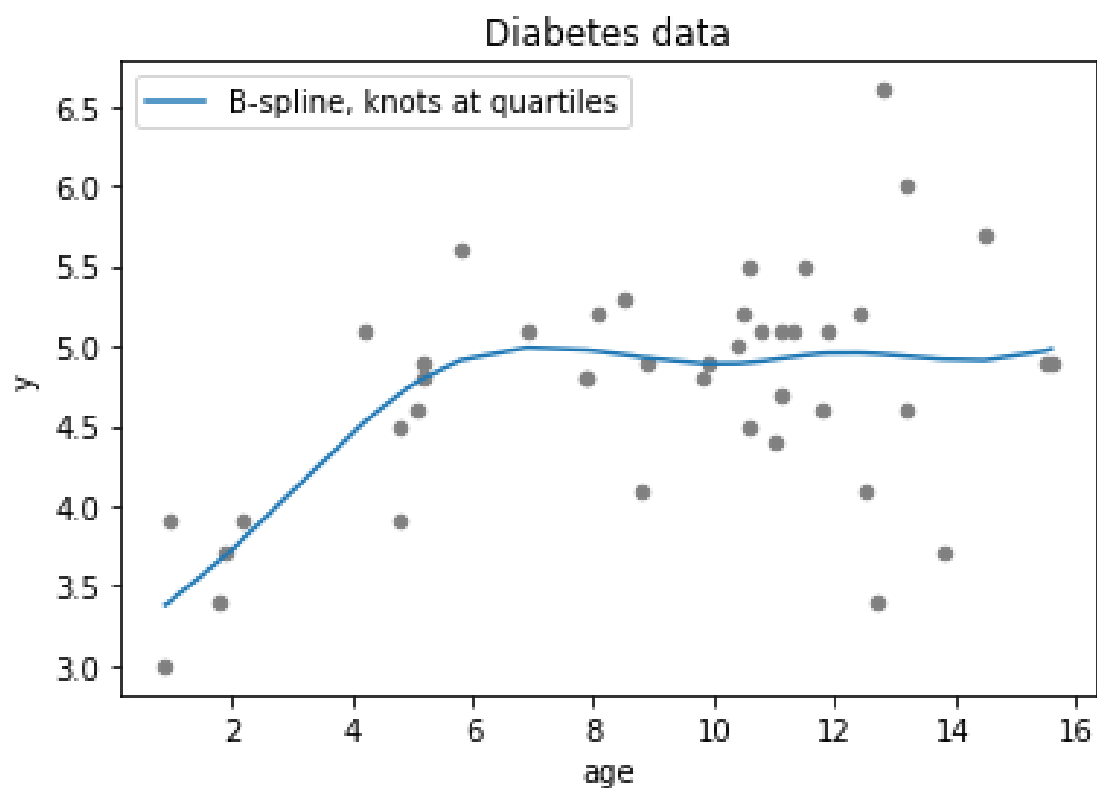
Related comments to cubic splines

- Instead of using the truncated cubic basis functions, it is more common in software packages to use a “B-spline” basis.
- This basis spans the same as the truncated cubic basis (and has $K+4$ basis functions), but is more efficient in estimation because at most five basis functions are involved in contributing towards any function value. The basis functions are defined recursively given the knots - difficult to write out.
- Because the behavior of cubic splines near boundaries can behave strangely, a variant is “natural” cubic splines. This basis enforces that the function is linear outside the extreme knots.

This linearity requirement creates 4 constraints (changing a cubic polynomial to linear at each extreme), so with K knots the total number of basis functions is K .

B-spline basis





Choosing the number and position of knots

- Can place knots at quantiles of the predictor variable, or at regularly spaced intervals
- Choosing the number of knots seems to be more crucial in obtaining a good non-linear fit
- Good idea to place more knots in parts of the predictor data where one might expect the mean function to change rapidly.

Regression splines and multiple predictors

Because a single predictor transformed as a polynomial term or a spline involves a linear combination of basis functions, it is actually straightforward to extend this approach to multiple predictors.

- Transform each predictor/feature individually to a linear combination of basis functions
- Include all terms simultaneously into the model
 - Fit least-squares multiple regression (if outcome is quantitative)
 - Fit multiple logistic regression (if outcome is binary)

Smoothers and additive models

So far, we have explored incorporating non-linear functions of predictors through polynomial terms or splines.

Notation

Let η_i denote the simultaneous contribution of the J predictors $x_{i1}, x_{i2}, \dots, x_{iJ}$ to the model.

For least-squares regression, we have $\eta_i = E(Y_i|x)$, and for logistic regression we have $\eta_i = \text{logit Pr}(Y_i = 1)$.

Example simultaneous contribution of features:

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_J x_{iJ}$$

for models where each predictor enters linearly.

Another example Let

$$\eta_i = \beta_0 + f_1(x_{i1}) + \dots + f_J(x_{iJ})$$

where $f_j(x_{ij})$ might be a cubic spline for variable x_j .

A model that includes the additive contribution of nonlinearly transformed predictors is called an additive model.

Actually, to be more precise, if the outcome variable is assumed to be normally distributed, then the above is an additive model. If the outcome variable is not normally distributed (e.g., binary), then the above model is called a generalized additive model.

Beyond polynomials and splines: Smoothers

Our approach thus far of identifying a non-linear function of a predictor has relied on

- Polynomial functions of the predictor
- Splines, i.e., piecewise polynomials that are connected at knots

Is there a more general way to view non-linear functions of a predictor?

Smoothers

A smoother $S(x)$ is a function that is used to estimate some feature of a response as a function of one or more predictors x that is typically less variable than the response.

By convention, we'll denote $s(x)$ to be an estimate of $S(x)$, called a “smooth,” based on observed data.

A “scatterplot smoother” is a smoother of just one predictor variable.

Example smoother: Lowess (or loess)

Abbreviation for “locally-weighted scatterplot smoother.”

To determine the smooth value at any value x ,

- Specify a kernel function, $K(x_i, x) \geq 0$. A kernel function is chosen so that
 - $K(x_i, x)$ is a maximum when $x = x_i$
 - $K(x_i, x)$ goes to 0 as $|x - x_i|$ increases
- At each x , perform a weighted least-squares regression of all points with kernel weights defined by $K(x_i, x)$.
- The smooth is determined as the fitted value of weighted least-squares regression at each location x .

A commonly used smoother: (Cubic) smoothing splines

Choice of a smoother as an optimization problem:

Among all functions f that are twice-differentiable, find the one that minimizes the penalized sum of squares

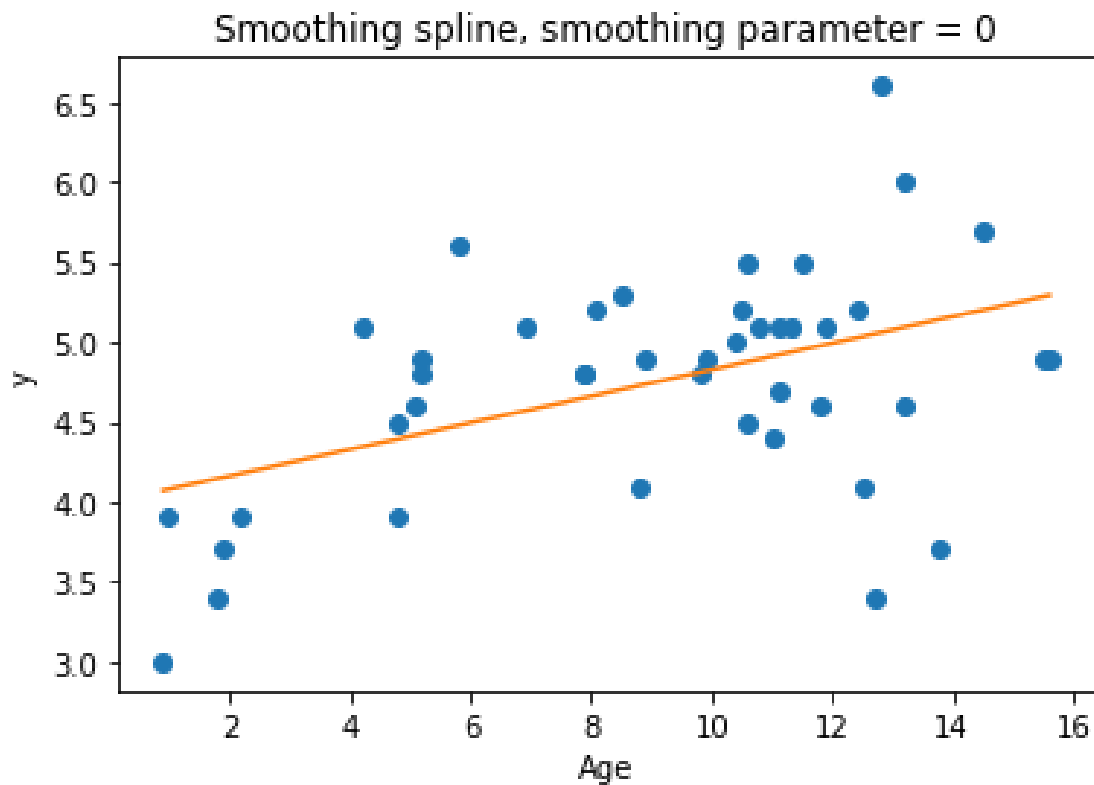
$$\text{PSS}(f \mid \mathbf{y}, \mathbf{x}) = p \sum_{i=1}^n (y_i - f(x_i))^2 + (1 - p) \int_a^b (f''(t))^2 dt$$

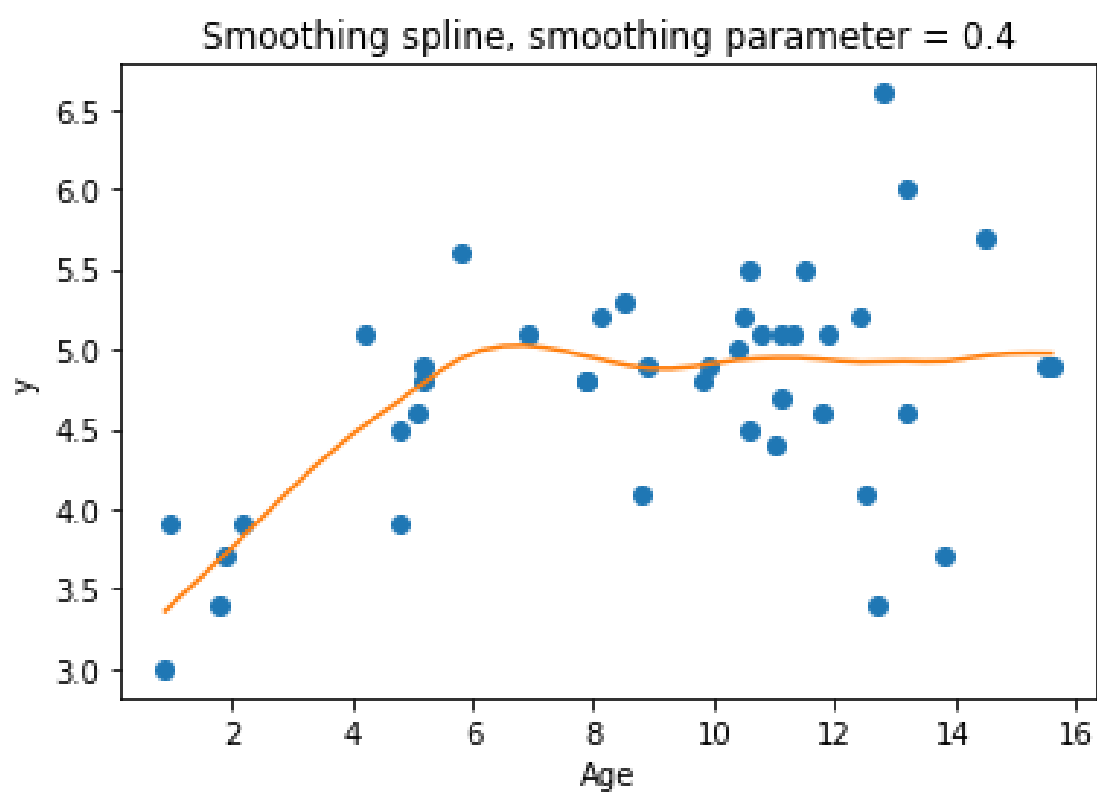
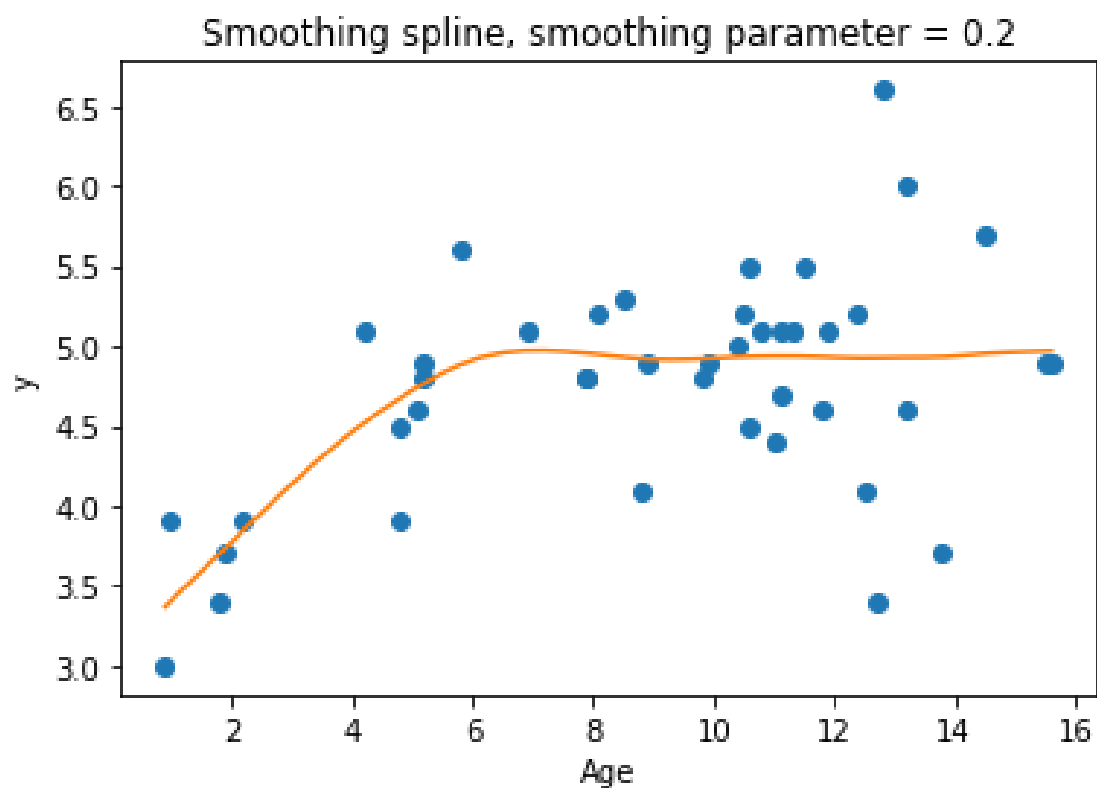
for a given value of $0 \leq p \leq 1$, where

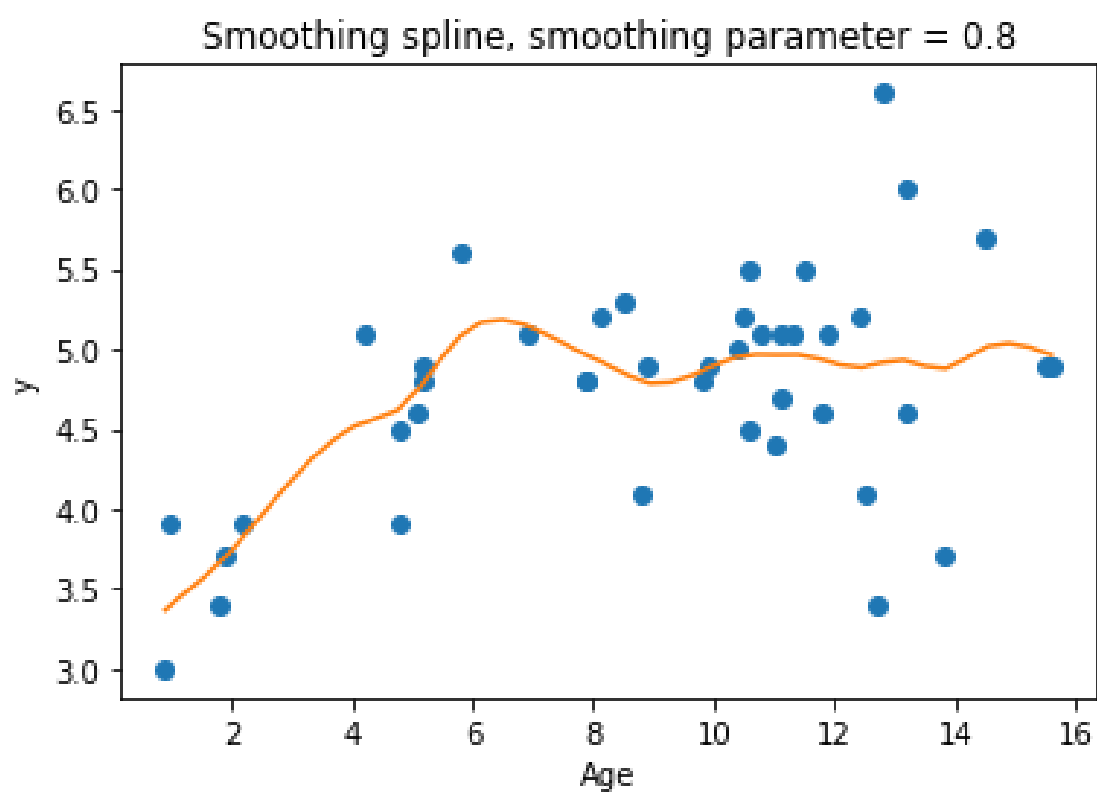
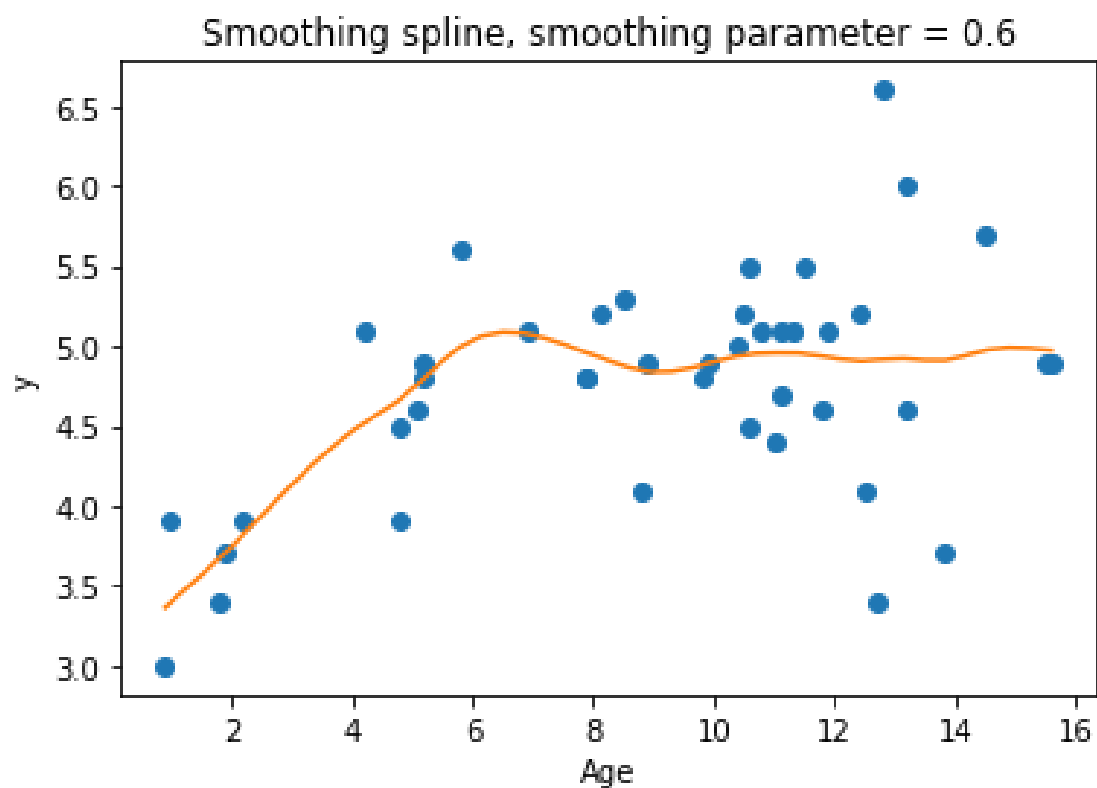
$$a \leq x_1 \leq x_2 \leq \cdots \leq x_n \leq b.$$

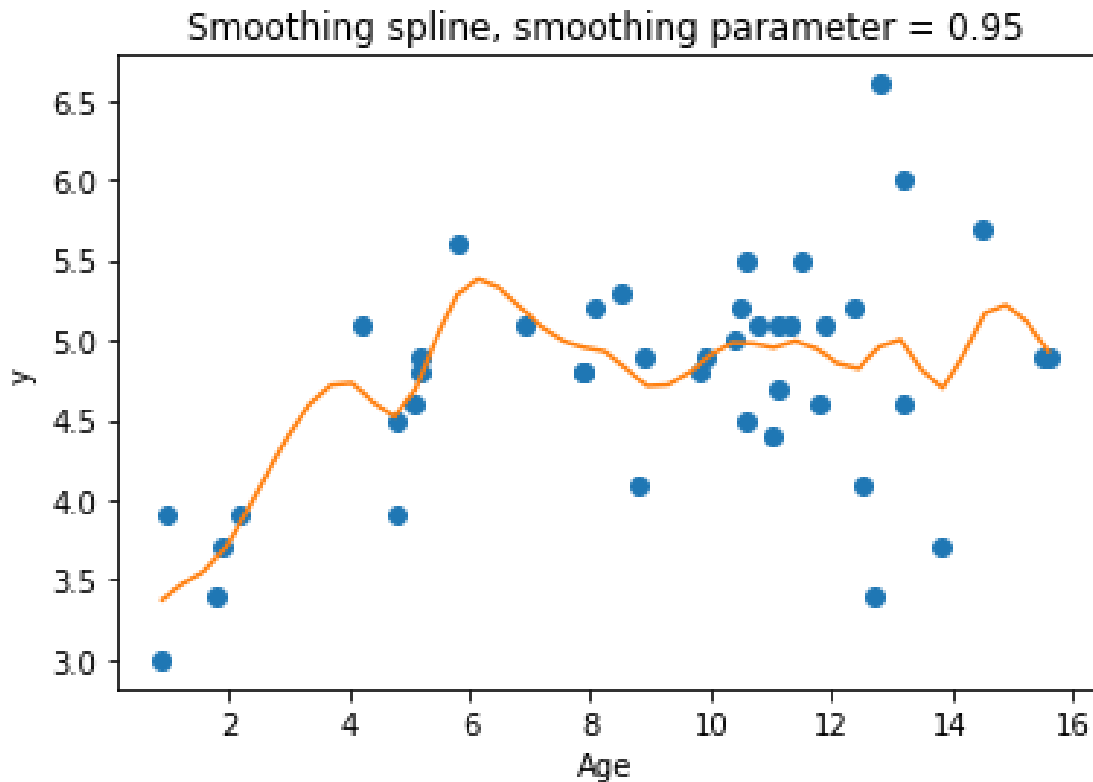
Why is this a reasonable criterion for f ?

- The first term measures the lack of fit having chosen f .
- The second term is a penalty for the wiggleness of f .
- The value p is a “tuning parameter” that balances the trade-off between lack of fit and wiggleness.
 - Large values of p result in a “connect-the-dots” fit.
 - Small values of p result in a fit closer to a line.









Solution to optimization problem

For a given s , the unique f that minimizes the penalized sum of squares turns out to be a cubic spline with knots at each x_i . As a reminder:

- A (different) cubic polynomial is fit between each x_i and x_{i+1} (not necessarily connecting to (x_i, y_i) or (x_{i+1}, y_{i+1})).
- The smoothing spline is continuous at each x_i .
- Not only do the cubic polynomials connect at each x_i , but the first and second derivatives of the cubic polynomials connecting at each x_i are equal. This is equivalent to saying that the smoothing spline is twice differentiable (even at the x_i).
- Conventionally we use assume a natural cubic spline: This smoothing spline is linear to the left of x_1 and to the right of x_n .

Actually, this is a special type of cubic spline called a “cubic smoothing spline.”

- Without further discussion, the cubic spline is over-parameterized. If the knots are at every x_i , there is no information to choose from among different cubic polynomials between x_i and x_{i+1} .
- The parameter p is what makes the connecting cubic polynomial choices unique.

Once p is specified, determining the smooth is just fitting a cubic spline.

Another possibility: B-splines and P-splines

- The B-spline basis of cubic polynomials is a perfectly reasonable approach to smoothing, as we have seen.
- Penalized B-splines, also known as P-splines, add a penalty function to reduce the impact of overfitting.
- The `PyGam` library, the main library for implementing GAMs, by default assumes
 - a cubic P-spline basis for constructing smoothers,
 - 20 functions in the P-spline basis (can be increased or decreased)

Unfortunately, `PyGam` does not implement smoothing splines.

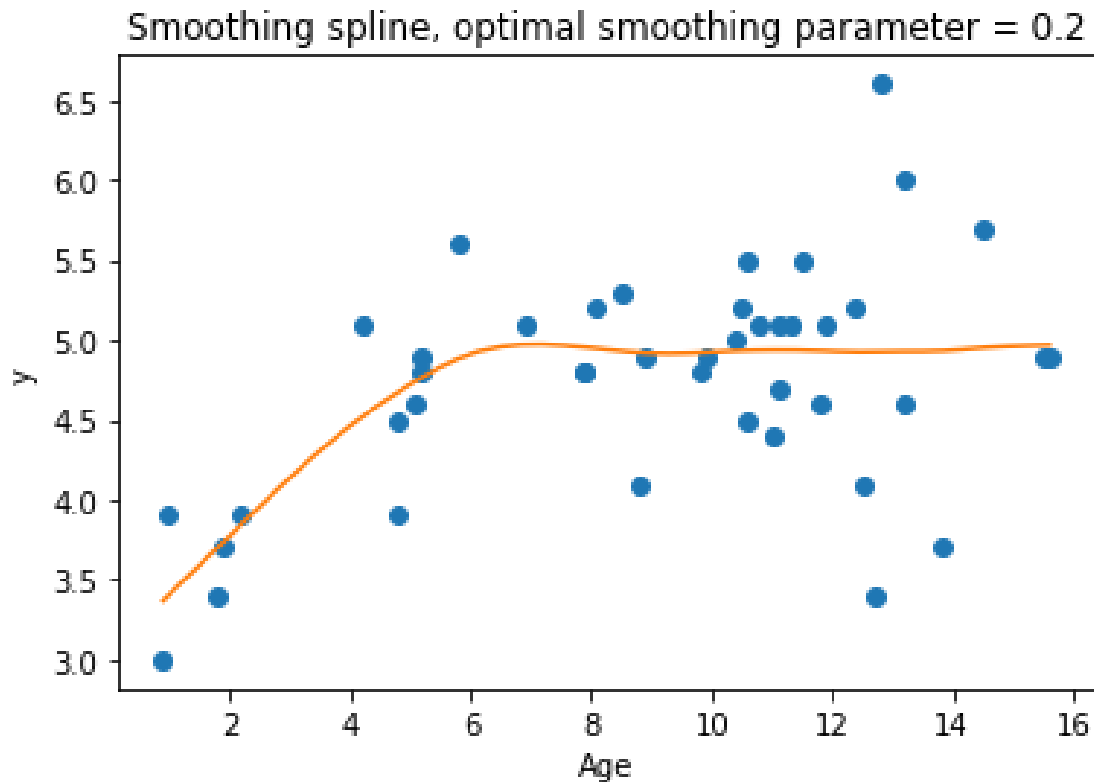
Choosing the smoothing parameter/penalty: Cross-validation (CV)

In general: Choose a smoothing parameter that produces the best predictions on data not analyzed.

For our setup when we observed pairs (x_i, y_i) , a reasonable CV criterion is

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - s_{\lambda}^{(i)}(x_i))^2$$

where λ is a tuning parameter, $s_{\lambda}^{(i)}$ is the smooth estimated from all the data not including observation i . Choose λ to minimize this expression.



Extending smoothers to multiple predictors

- Additive models: $y_i = \beta_0 + S_1(x_{i1}) + \dots + S_J(x_{iJ}) + \varepsilon_i$
- Generalized additive models (GAMs): Assume the Y_i have a distribution such as binomial, with $\mu_i = E(Y_i|x)$, with the additive predictor

$$\text{logit Pr}(Y_i = 1) = \eta_i = \beta_0 + S_1(x_{i1}) + \dots + S_J(x_{iJ})$$

Additive models are a special case of GAMs. Also, logistic regression is a special case of GAMs.

Worth noting that GAMs can include linear terms like $x_j\beta_j$ in addition to the smoother terms ($x_j\beta_j$ is a particular type of smoother).

Fitting a GAM

Start by picking arbitrary $\hat{\beta}_0, s_1(x_1), \dots, s_J(x_J)$. From this information, we can compute the $\hat{\eta}_i, \hat{\mu}_i$, etc.

- Define the i -th working response

$$z_i = \begin{cases} \hat{\eta}_i + (y_i - \hat{\mu}_i)(\hat{\mu}_i(1 - \hat{\mu}_i)) & \text{for logistic regression models} \\ y_i & \text{for least-squares regression models} \end{cases}$$

- For fixed $\mathbf{z} = (z_1, \dots, z_n)$, run “backfitting algorithm.” That is, smooth

$$\mathbf{z} - \hat{\beta}_0 - s_1(\mathbf{x}_1) - \dots - s_{j-1}(\mathbf{x}_{j-1}) - s_{j+1}(\mathbf{x}_{j+1}) - \dots - s_J(\mathbf{x}_J)$$

on \mathbf{x}_j to get new estimated $s_j(\mathbf{x}_j)$. Iterate 3-4 times to get improved estimates of the $s_j(\mathbf{x}_j)$ for all j .

Repeat the above two steps until convergence.

Cross-validation criterion for smoothing parameter for binary response models

Instead of choosing λ to minimize

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n (y_i - s_{\lambda}^{(i)}(x_i))^2,$$

minimize the sum

$$CV(\lambda) = - \sum_{i=1}^n \left(y_i \log p_{\lambda}^{(i)} + (1 - y_i) \log(1 - p_{\lambda}^{(i)}) \right).$$

Example: Kyphosis data

Data set consists of 81 children who have had corrective spinal surgery.

The variables:

- **Kyphosis:** a binary factor indicating if a kyphosis (type of deformation) was present after the operation
- **Age:** age in months
- **Number:** the number of vertebrae involved
- **Start:** the number of the topmost vertebra operated on

Data summaries

	Age	Number	Start
count	81.000000	81.000000	81.000000
mean	83.654321	4.049383	11.493827
std	58.104251	1.619423	4.883962
min	1.000000	2.000000	1.000000
25%	26.000000	3.000000	9.000000
50%	87.000000	4.000000	13.000000
75%	130.000000	5.000000	16.000000
max	206.000000	10.000000	18.000000

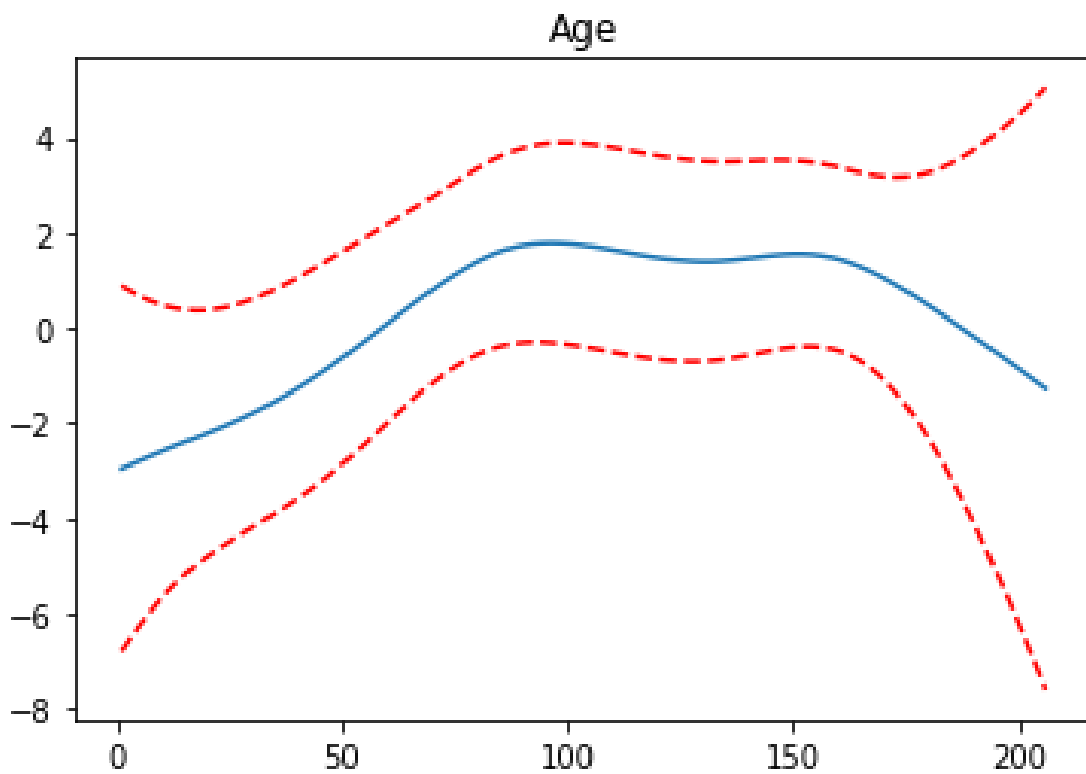
GAM in action

```

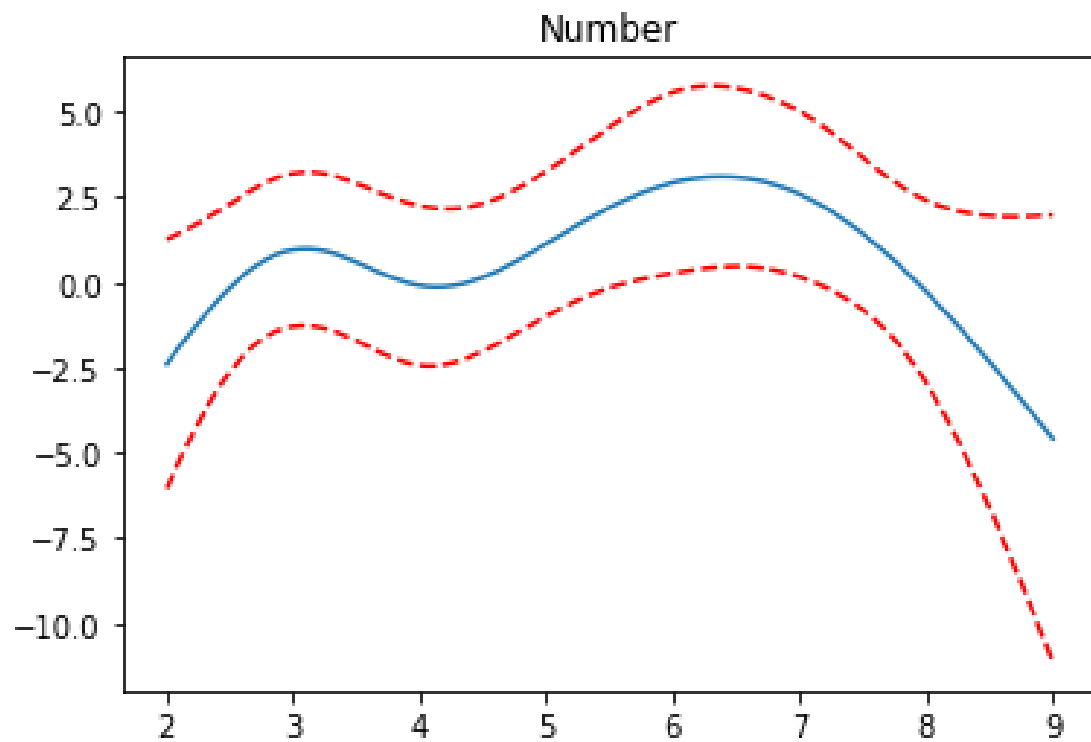
LogisticGAM
=====
Distribution:          BinomialDist Effective DoF:          13.4257
Link Function:         LogitLink   Log Likelihood:        -13.3818
Number of Samples:     64 AIC:          53.6151
                        AICc:         62.7774
                        UBRE:         3.0056
                        Scale:         1.0
                        Pseudo R-Squared: 0.5857
=====
Feature Function      Lambda      Rank      EDoF      P > x      Sig. Code
=====
s(0)                  [0.6]      20       7.5      8.73e-01
s(1)                  [0.5]      20       4.8      5.99e-02    .
s(2)                  [0.6]      20       1.2      6.98e-02    .
intercept             1         1        0.0      9.90e-03    **
=====
Significance codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

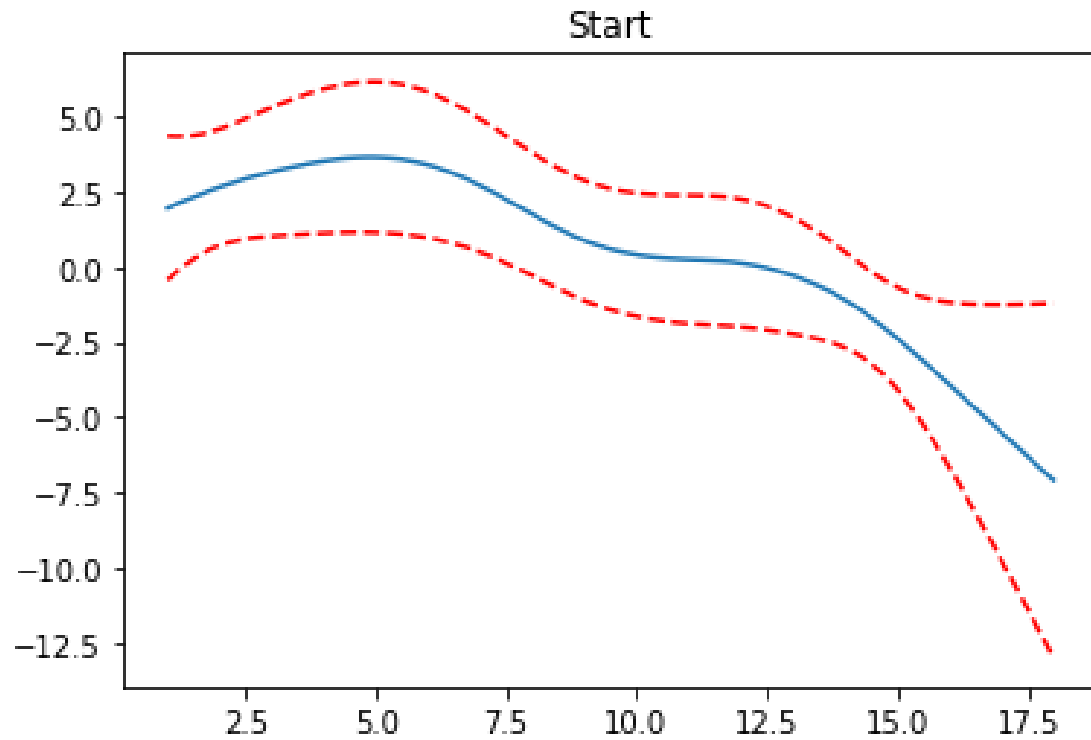
Kyphosis: Smooth term Age



Kyphosis: Smooth term Number



Kyphosis: Smooth term Start



Prediction from GAM

Can apply the fit of a GAM to make prediction estimates on new data.

Simple experiment Want to know which of the following two GAMs is preferable:

- GAM including `Age`, `Number` and `Start` as features, or
- GAM including only `Age` and `Start` as features.

Strategy

- Split observations into training and test data (e.g., 80% training set).
Stratify based on the binary outcome to ensure balance.
- Fit each GAM to the training data, and make probability predictions to the test data.
- Compute correct classification rate for each GAM as the approach to compare prediction accuracy. (Other approaches possible as well)

Test accuracy

- Three-feature model: 82% correct classification
- Two-feature model: 76% correct classification

However, with such a small test set (16 observations) the difference is not conclusive.

See accompanying python code for evaluating via cross-validation which also gives inconclusive results.