

# Chapter 7: Moving Beyond Linearity

So far we have mainly focused on linear models.

Linear models are relatively simple to describe and implement.

+ : interpret + inference

- : can have limited predictive performance because linearity assumption is always an approximation (may not be a good one).

Previously, we have seen we can improve upon least squares using ridge regression, the lasso, principal components regression, and more.

improvement obtained by reducing complexity of linear models  $\Rightarrow$  lower variance of estimates.  
still a linear model! can only be improved so much.

Through simple and more sophisticated extensions of the linear model, we can relax the linearity assumption while still maintaining as much interpretability as possible.

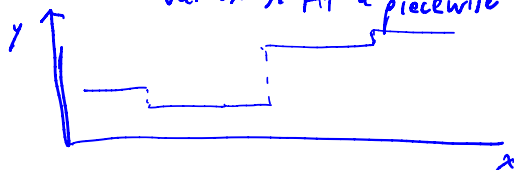
① Polynomial regression: adding extra predictors that are original variables raised to a power

e.g. cubic regression use  $X, X^2, X^3$  as predictors, e.g.  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon$

+ : non-linear fit

- : with large powers, polynomial can take very strange shapes (especially at boundary).

② Step functions: cut the range of predictor into  $K$  distinct regions (to produce categorical variable). Fit a piecewise constant function to (binned)  $X$ .



③ Regression Splines: more flexible than polynomials + step functions (extends both)

idea: cut range of  $X$  into  $K$  distinct regions + polynomial is fit within each region  
polynomials constrained so they smoothly joined.

④ Generalized additive models: extend above ideas to deal w/ multiple predictors.

Note: We can talk about regression or classification, e.g. logistic regression (polynomial):  $P(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 X + \dots + \beta_d X^d)}{1 + \exp(\beta_0 + \beta_1 X + \dots + \beta_d X^d)}$

We've already seen this.

# 1 Step Functions

Using polynomial functions of the features as predictors imposes a global structure on the non-linear function of  $X$ .

We can instead use *step-functions* to avoid imposing a global structure.

idea: Break range of  $X$  into bins and fit different constant to each bin.

details: ① create cut points  $c_1, \dots, c_k$  in the range of  $X$

② Construct  $k+1$  new variables.

$$c_0(x) = \mathbb{I}(X < c_1)$$

$$c_1(x) = \mathbb{I}(c_1 \leq X < c_2)$$

⋮

$$c_{k-1}(x) = \mathbb{I}(c_{k-1} \leq X < c_k)$$

$$c_k(x) = \mathbb{I}(c_k \leq X)$$

indicator variable  
 $c_0(x) = \mathbb{I}(X < c_1) = \begin{cases} 1 & X < c_1 \\ 0 & \text{o.w.} \end{cases}$   
"dummy variables"

③ Use least squares to fit a linear model using  $c_1(x), c_2(x), \dots, c_k(x)$

$$Y = \beta_0 + \beta_1 c_1(x) + \dots + \beta_k c_k(x) + \varepsilon$$

↑ leave out  $c_0(x)$  because it is equivalent to including an intercept.

For a given value of  $X$ , at most one of  $c_1, \dots, c_k$  can be non-zero.

$c_0(x) + c_1(x) + \dots + c_k(x) = 1$  since  $X$  must be in exactly one interval.

When  $X < c_1 \Rightarrow$  all of predictors  $c_1, \dots, c_k = 0$

$\Rightarrow \beta_0$  interpreted as the mean value of  $Y$  when  $X < c_1$

$\beta_j$  represent the average increase in the response for  $c_j \leq X < c_{j+1}$  relative to  $X < c_1$ .

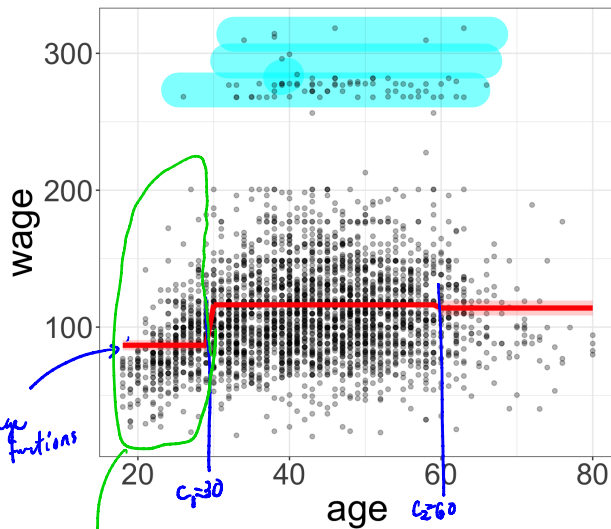
We can also fit a logistic regression model for classification.

$$P(Y=1|X) = \frac{\exp(\beta_0 + \beta_1 c_1(x) + \dots + \beta_k c_k(x))}{1 + \exp(\beta_0 + \beta_1 c_1(x) + \dots + \beta_k c_k(x))}$$

Example: Wage data. *for a group of 3000 male workers in mid-atlantic region*

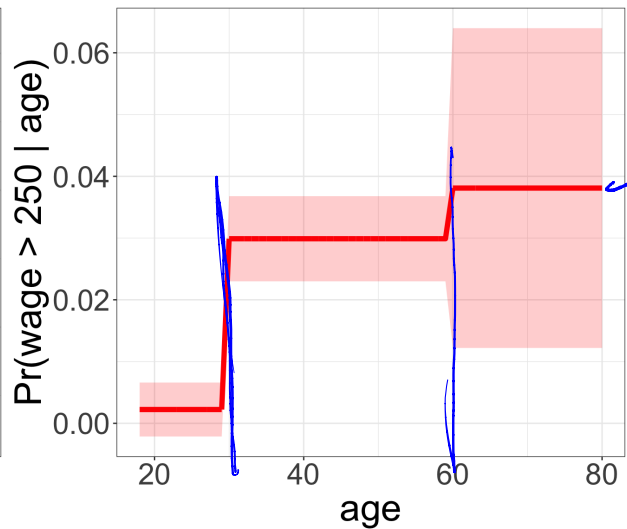
year	age	maritl	race	edu- cation	region	job- class	health	health_ins	logwage	wage
2006	18	1. Never Married	1. White	1. < HS Grad	2. Middle Atlantic	1. Industrial	1. <=Good	2. No	4.318063	75.04315
2004	24	1. Never Married	1. White	4. College Grad	2. Middle Atlantic	2. Information	2. >=Very Good	2. No	4.255273	70.47602
2003	45	2. Married	1. White	3. Some College	2. Middle Atlantic	1. Industrial	1. <=Good	1. Yes	4.875061	130.98218
2003	43	2. Married	3. Asian	4. College Grad	2. Middle Atlantic	2. Information	2. >=Very Good	1. Yes	5.041393	154.68529

$c_1 = 30$   
 $c_2 = 60$



*filled value of wage using step functions of age.*

*missing clear upward trend.*



*logistic regression modeling probability of being high wage earner given age. stepwise fit model w/ knots at X=30, 60.*

*Unless there are natural breakpoints in the predictor, piecewise constant can miss trends.*

## 2 Basis Functions

Polynomial and piecewise-constant regression models are in fact special cases of a *basis function approach*.

**Idea:**

have a family of functions or transformations that can be applied to a variable  $X$   
 $b_1(x), b_2(x), \dots, b_k(x)$

Instead of fitting the linear model in  $X$ , we fit the model

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \dots + \beta_k b_k(x_i) + \varepsilon_i$$

Note that the basis functions are fixed and known. we choose them ahead of time.

ex: polynomial regression  $b_j(x_i) = x_i^j, j=1, \dots, d.$

ex: step function  
(piecewise constant)

$$b_j(x_i) = \mathbb{I}(c_j \leq x_i < c_{j+1}) \\ = \begin{cases} 1 & c_j \leq x_i < c_{j+1} \\ 0 & \text{o.w.} \end{cases}$$

We can think of this model as a standard linear model with predictors defined by the basis functions and use least squares to estimate the unknown regression coefficients.

$\Rightarrow$  we can use all our inference tools for linear models.

e.g.  $se(\hat{\beta}_j)$  and  $F$ -statistic for model significance.

Many alternatives exist for basis functions.

e.g. wavelets, fourier series, regression splines (next).

# 3 Regression Splines

Regression splines are a very common choice for basis function because they are quite flexible, but still interpretable. Regression splines extend upon polynomial regression and piecewise constant approaches seen previously.

start

## 3.1 Piecewise Polynomials

Instead of fitting a high degree polynomial over the entire range of  $X$ , piecewise polynomial regression involves fitting separate low-degree polynomials over different regions of  $X$ .

e.g. one knot at  $c$

fit two polynomials to the data  
one on subset for  $x < c$   
one on subset for  $x \geq c$

each polynomial can be fit using least squares.

For example, a piecewise cubic with no knots is just a standard cubic polynomial.

A piecewise cubic with a single knot at point  $c$  takes the form

$$y_i = \begin{cases} \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \epsilon_i & \text{if } x_i < c \\ \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \epsilon_i & \text{if } x_i \geq c \end{cases}$$

Using more knots leads to a more flexible piecewise polynomial.

if we place  $k$  knots  $\Rightarrow$  fit  $k+1$  polynomials

In general, we place  $K$  knots throughout the range of  $X$  and fit  $K + 1$  polynomial regression models.

This leads to  $(d+1)(k+1)$  degrees of freedom in model  
(# of parameters to fit  $\approx$  complexity / flexibility).

## 3.2 Constraints and Splines

To avoid having too much flexibility, we can *constrain* the piecewise polynomial so that the fitted curve must be continuous.

*i.e. there cannot be a jump at knots.*

To go further, we could add two more constraints

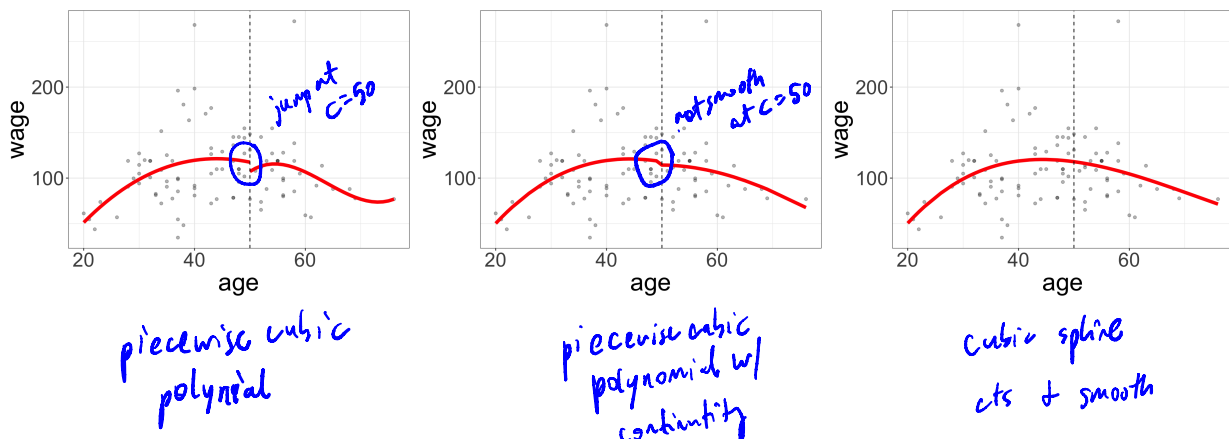
- ① 1st derivative of piecewise polynomial must be continuous
- ② 2nd derivative of piecewise polynomial must be cts.

In other words, we are requiring the piecewise polynomials to be smooth.

Each constraint that we impose on the piecewise cubic polynomials effectively frees up one degree of freedom, by reducing the complexity of the resulting fit.

<sup>cubic</sup>  
The fit with continuity and 2 smoothness constraints is called a <sup>cubic</sup> *spline*.

A degree- $d$  spline is a piecewise degree- $d$  polynomial w/ continuity in derivatives up to degree  $d-1$  at each knot.



### 3.3 Spline Basis Representation

Fitting the spline regression model is more complex than the piecewise polynomial regression. We need to fit a degree  $d$  piecewise polynomial and also constrain it and its  $d - 1$  derivatives to be continuous at the knots.

We can use the basis model to represent a regression spline.

eg. cubic spline w/  $K$  knots

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K+3} b_{K+3}(x_i) + \epsilon_i$$

for appropriate basis functions  $b_1, b_2, \dots, b_{K+3}$

$x, x^2, x^3$

$\bullet = d$ -dimensional spline.

basis for  $d$  degree polynomial

The most direct way to represent a cubic spline is to start with the basis for a cubic polynomial and add one truncated power basis function per knot.

$$h(x, \xi) = (x - \xi)_+^3 = \begin{cases} (x - \xi)^3 & \text{if } x > \xi \\ 0 & \text{o.w.} \end{cases} \quad \text{where } \xi \text{ is a knot.}$$

$$\Rightarrow y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \sum_{j=1}^K \beta_{3+j} h(x_i, \xi_j)$$

this will lead to discontinuity in only the 3rd derivative at each  $\xi_j$  with continuous first and second derivatives and continuity at each  $\xi_j$

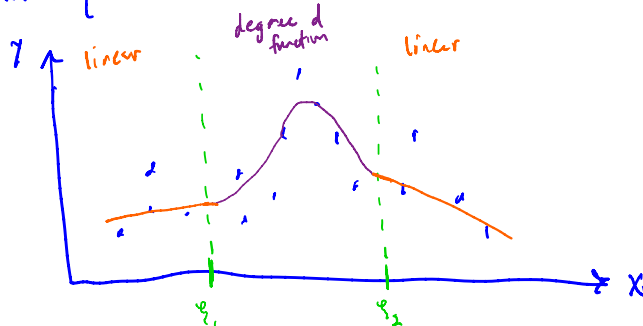
$$df = K + 4 \quad (\text{cubic spline w/ } K \text{ knots})$$

Unfortunately, splines can have high variance at the outer range of the predictors. One solution is to add boundary constraints.

when  $X$  is very small or very large.

$\Rightarrow$  "natural spline"

function required to be linear at the boundary (where  $X$  is smaller than the smallest knot and bigger than biggest knot)



additional constraint produces more stable estimates at the boundaries.

### 3.4 Choosing the Knots

When we fit a spline, where should we place the knots?

⇒ regression spline is most flexible in regions that contain a lot of knots (coefficients changing rapidly).  
 ⇒ place knots where we think the function will vary rapidly and less knots where function is stable.

More common in practice: place them uniformly.

To place knots: choose desired degrees of freedom (flexibility) & use software to automatically place # knots at uniform quantiles of data.

How many knots should we use?

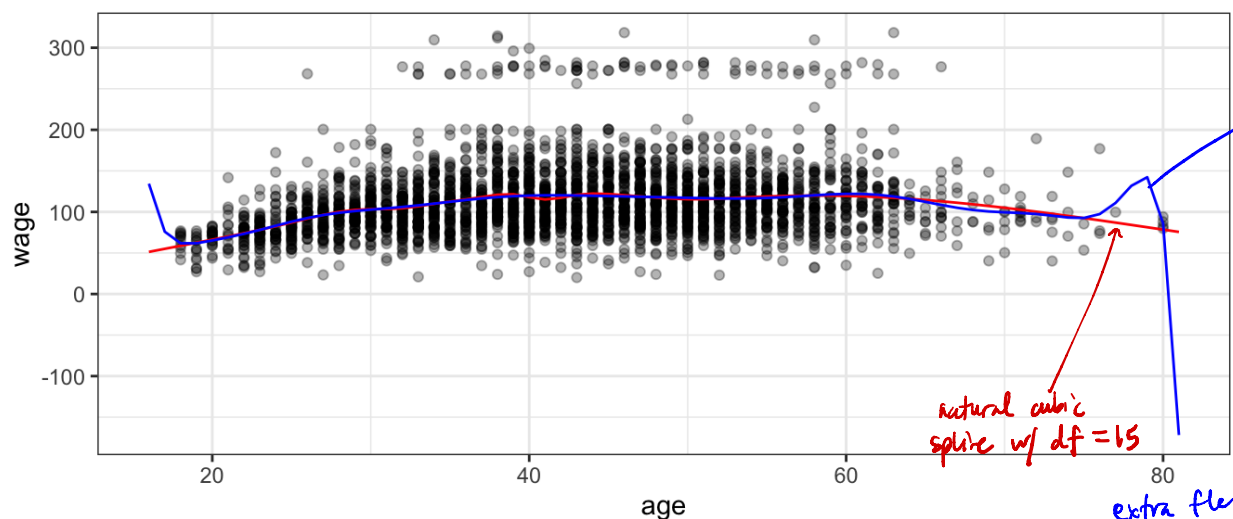
⇔ how many df should we use?

Use CV! Use  $k$  giving smallest CV MSE (CV error).

### 3.5 Comparison to Polynomial Regression

Regression splines often give superior results to polynomial regression.

↳ Polynomial regression must use high degrees to achieve flexibility (e.g.  $X^{15}$ ), but regression splines introduce flexibility through knots (fixed degree polynomials) ⇒ more stability (esp. at boundaries)



extra flexibility of polynomial at boundary produces undesirable results, but NC spline w/ same flexibility still looks reasonable.



# 4 Generalized Additive Models

So far we have talked about flexible ways to predict  $Y$  based on a single predictor  $X$ .

These approaches can be seen as extensions of simple linear regression model.

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

extension: basis functions of  $X$

*Generalized Additive Models (GAMs)* provide a general framework for extending a standard linear regression model by allowing non-linear functions of each of the variables while maintaining *additivity*.

flexibly predict  $Y$  on several predictors  $X_1, \dots, X_p$ .

## 4.1 GAMs for Regression

still additive models

can be used for regression or classification.

A natural way to extend the multiple linear regression model to allow for non-linear relationships between feature and response:

linear regression:  $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i$

extension idea: non-linear replace each linear component  $\beta_j x_{ij}$  with a smooth non-linear function.

$$\Rightarrow \text{GAM: } y_i = \beta_0 + \sum_{j=1}^p f_j(x_{ij}) + \varepsilon_i$$

$$\approx \beta_0 + f_1(x_{i1}) + f_2(x_{i2}) + \dots + f_p(x_{ip}) + \varepsilon_i$$

"additive" because we calculate a separate  $f_j$  for each predictor  $X_j$  and add them together.

possibilities for  $f_j$ :

- linear component (leads to linear regression).

- polynomial function

- regression spline

- smoothing spline

- local linear regression.

] not covered, but see textbook Ch. 7.5-7.6 for details.

The beauty of GAMs is that we can use our fitting ideas in this chapter as building blocks for fitting an additive model.

Example: Consider the Wage data.

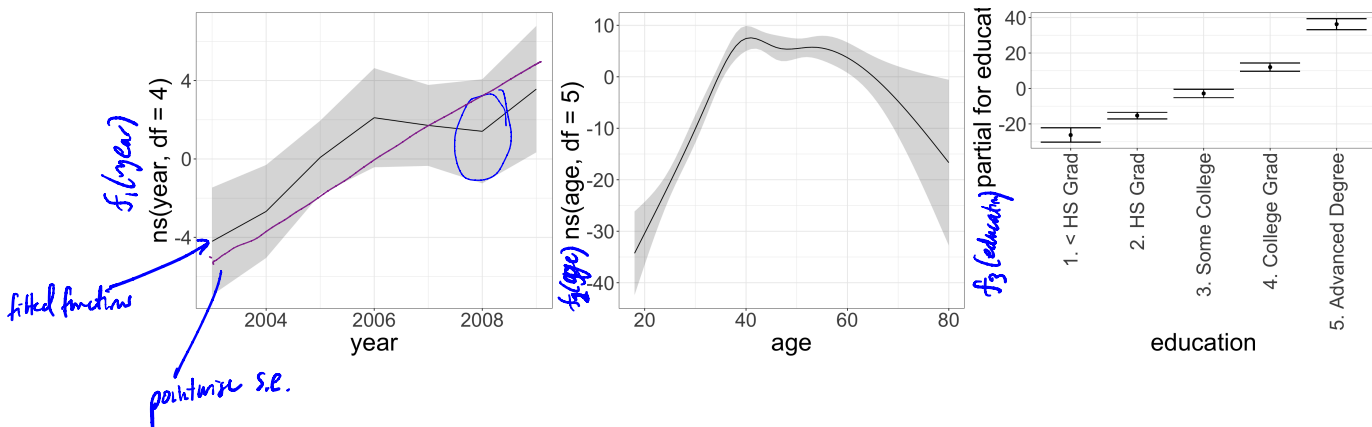
$$\text{Wage} = \beta_0 + f_1(\text{year}) + f_2(\text{age}) + f_3(\text{education}) + \varepsilon$$

← quantitative
← categorical.

where  $f_1$  is natural spline w/ 4 df

$f_2$  is natural cubic spline w/ 5 df

$f_3$  constant functions for each value (dummy variables)



relationship between each variable and response.

- year: holding age and education fixed, wage tends to increase with year. (inflation?).
- age: holding year and education fixed, wage is low for young and old people, highest for intermediate ages.
- education: holding year & age fixed, wage tends to increase with education.

We could easily replace  $f_j$  with different smooth functions and get different fits.

just need to change basis & use least squares.

## Pros and Cons of GAMs

Advantages

- allow nonlinear fits  $f_j$  for each  $X_j$  to automatically model non-linear relationship that linear regression will miss.
  - non-linear fit can lead to more accurate prediction of response (if there is a truly non-linear relationship)
  - additive model  $\Rightarrow$  we can still examine effect of each  $X_j$  on response individually holding all others fixed.
- $\Rightarrow$  GAMs provide a useful representation for inference/interpretation.
- smoothness of  $f_j$  can be summarized by df.

Limitations

- model is restricted to be additive  
i.e. we can miss important interactions

solution: as with linear regression, we can manually add interactions by including additional predictors of the form  $X_j \cdot X_k$   
or add low-dim interaction terms of the form  $f_{jk}(X_j, X_k)$ .

$\nearrow$  two-dimensional splines (not covered).

For fully general models, we have to look for even more flexible approaches like random forests or boosting (next week!).

GAMs provide a useful compromise between linear and nonparametric approaches.

## 4.2 GAMs for Classification

assume  $Y$  takes values 0 or 1 (generalizations exist)

GAMs can also be used in situations where  $Y$  is categorical. Recall the logistic regression model:

$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

"logit" or log odds of  $P(Y=1|X)$  vs.  $P(Y=0|X)$ .  
as linear function of predictors.

A natural way to extend this model is for non-linear relationships to be used.

$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + f_1(x_1) + \dots + f_p(x_p)$$

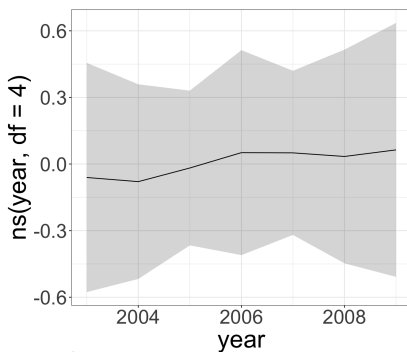
logistic regression GAM

Example: Consider the Wage data.

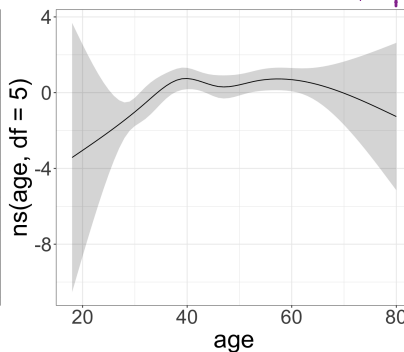
let  $Y = \text{wage} > \$250k$

we could fit a GAM:

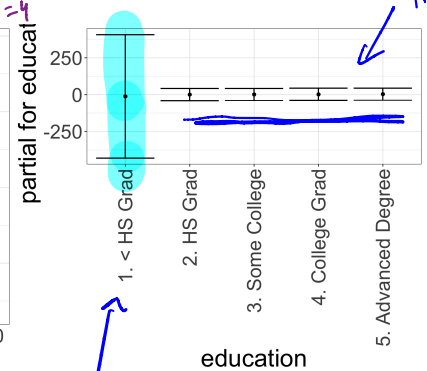
$$\log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + f_1(\text{year}) + f_2(\text{age}) + f_3(\text{education})$$



this looks quite linear could replace w/ linear function (polynomial w/ degree 1) without much  $\uparrow$  bias &  $\downarrow$  variance.



natural cubic spline w/ df=4



piecewise constants for each level. increase w/ education

based on scales of  $f_1, f_2, f_3$ , age & education have more of an effect on  $P(\text{higherner}(X))$  than year.

Nobody in the data set w/  $<$  HS education and wage  $>$  250k may want to refit this model excluding that class.