# Assignment 08 - STAT 689

*Philip Anderson; panders2@tamu.edu*

*2/26/2018*

```r
library("mgcv")
```

## Part 1

```r
# read in the modified salinity dataset
salinity <- read.csv('/Users/panders2/Documents/schools/tamu/stat_689/homework/semiparametric-regression
names(salinity) <- tolower(names(salinity))
str(salinity)
```
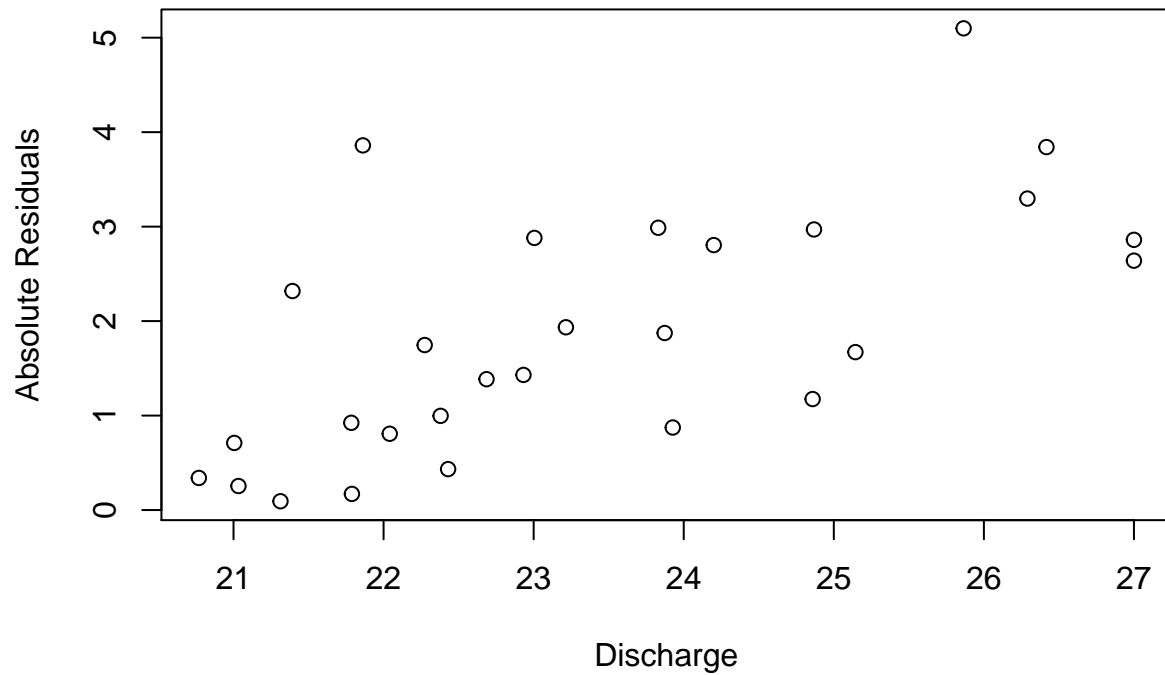
```
## 'data.frame':    28 obs. of  4 variables:
##  $ salinity       : num  7.6 7.7 4.3 5.9 5 6.5 8.3 8.2 13.2 12.6 ...
##  $ lagged.salinity: num  8.2 7.6 4.6 4.3 5.9 5 6.5 8.3 10.1 13.2 ...
##  $ trend          : int  4 5 0 1 2 3 4 5 0 1 ...
##  $ discharge      : num  23 23.9 26.4 24.9 27 ...
```

## Question 1

Fit a GAM where salinity is regressed on a smooth function of discharge. Plot discharge against the absolute residuals from this fit.

```r
gam_mod <- mgcv::gam(salinity ~ s(discharge)
                     , data=salinity
                     )
abs_resid <- abs(gam_mod$residuals)
plot(salinity$discharge, abs_resid
     , main="Discharge by Abs(Residual) from fitted GAM"
     , xlab="Discharge"
     , ylab="Absolute Residuals"
     )
```

## Discharge by Abs(Residual) from fitted GAM



Note that we appear to have a positive trend. This reflects evidence of hetereoscedacticity.

# Question 2

Calculate the ratio of the maximum predicted absolute residual to the minimum predicted absolute residual.

```
cat("\nMaximum Value: ", max(abs_resid))
```

```
##
## Maximum Value:  5.098749
```

```
cat("\nMinimum Value: ", min(abs_resid))
```

```
##
## Minimum Value:  0.09310846
```

```
cat("\nRatio of max to min: ", max(abs_resid) / min(abs_resid))
```

```
##
## Ratio of max to min:  54.7614
```

# Question 3

The ratio found in Question 2 is quite high (definitely higher than our heuristic of 3), indicating that we have evidence of heteroscedacticity. This will make the calculation of pointwise confidence intervals more involved, and render naive estimates incorrect.
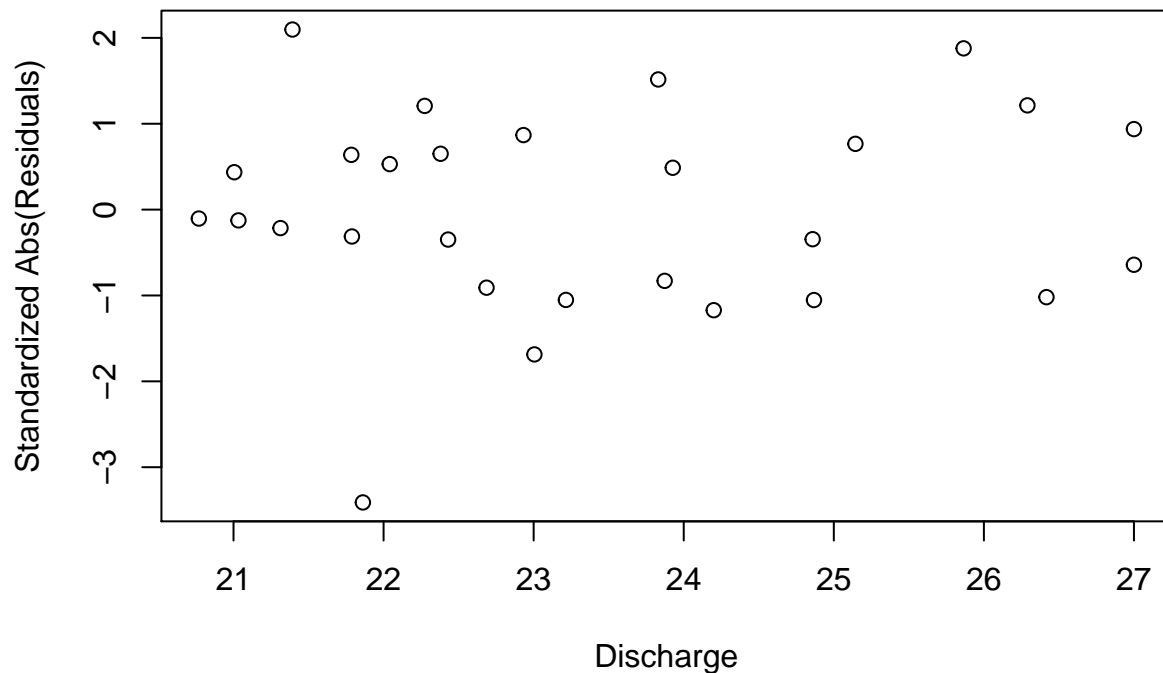
# Question 4

Conduct a weighted fit of the data. Plot discharge against the absolute residuals of the new model divided by the estimated standard deviations.

```
gam_mod_resid <- mgcv::gam(abs_resid ~ s(salinity$discharge))
weight <- 1/(fitted(gam_mod_resid)^2)
gam_mod_weight <- mgcv::gam(salinity ~ s(discharge)
                            , data=salinity
                            , weight=weight
                            )
```

```
plot(salinity$discharge, gam_mod_weight$residuals/fitted(gam_mod_resid)
     , main="Discharge against Standardized Absolute Residuals (Weighted Fit)"
     , xlab="Discharge"
     , ylab="Standardized Abs(Residuals)"
     )
```

**Discharge against Standardized Absolute Residuals (Weighted Fit)**



# Question 5

Compute the 95% confidence intervals for the unweighted and weighted fits.

```
x_var <- seq(from=min(salinity$discharge), to=max(salinity$discharge)
             , len=1000)
# predictions
unweight_pred <- predict(gam_mod
                         , newdata=data.frame(discharge=x_var)
                         , se.fit=TRUE)
```

3

```r
weight_pred <- predict(gam_mod_weight
                       , newdata=data.frame(discharge=x_var)
                       , se.fit=TRUE)

# confidence bounds
unweight_lower <- unweight_pred$fit - (1.96 * unweight_pred$se.fit)
unweight_upper <- unweight_pred$fit + (1.96 * unweight_pred$se.fit)

weight_lower <- weight_pred$fit - (1.96 * weight_pred$se.fit)
weight_upper <- weight_pred$fit + (1.96 * weight_pred$se.fit)
```

```r
plot(salinity$discharge, salinity$salinity, type="n"
     , main="Salinity by Discharge"
     , xlab="Discharge"
     , ylab="Salinity"
     )

gam_plotter <- function(pred_obj, lower_bnd, upper_bnd, shader) {
  lines(x_var, pred_obj$fit, lwd=2.5)
  polygon(x=c(x_var, rev(x_var))
          , y=c(upper_bnd, rev(lower_bnd))
          , col=shader
          , border=NA
          , density=50
          )
}

# unweighted GAM first
gam_plotter(pred_obj=unweight_pred
            , lower_bnd=unweight_lower
            , upper_bnd=unweight_upper
            , shader="palegreen"
            )
rug(salinity$discharge, lwd=3, col="olivedrab")

# weighted GAM second
gam_plotter(pred_obj=weight_pred
            , lower_bnd=weight_lower
            , upper_bnd=weight_upper
            , shader="cornflowerblue"
            )
legend("topright"
       , c("Unweighted", "Weighted")
       , lty=c(1,1)
       , lwd=c(5,5)
       , col=c("palegreen", "cornflowerblue")
       )
```
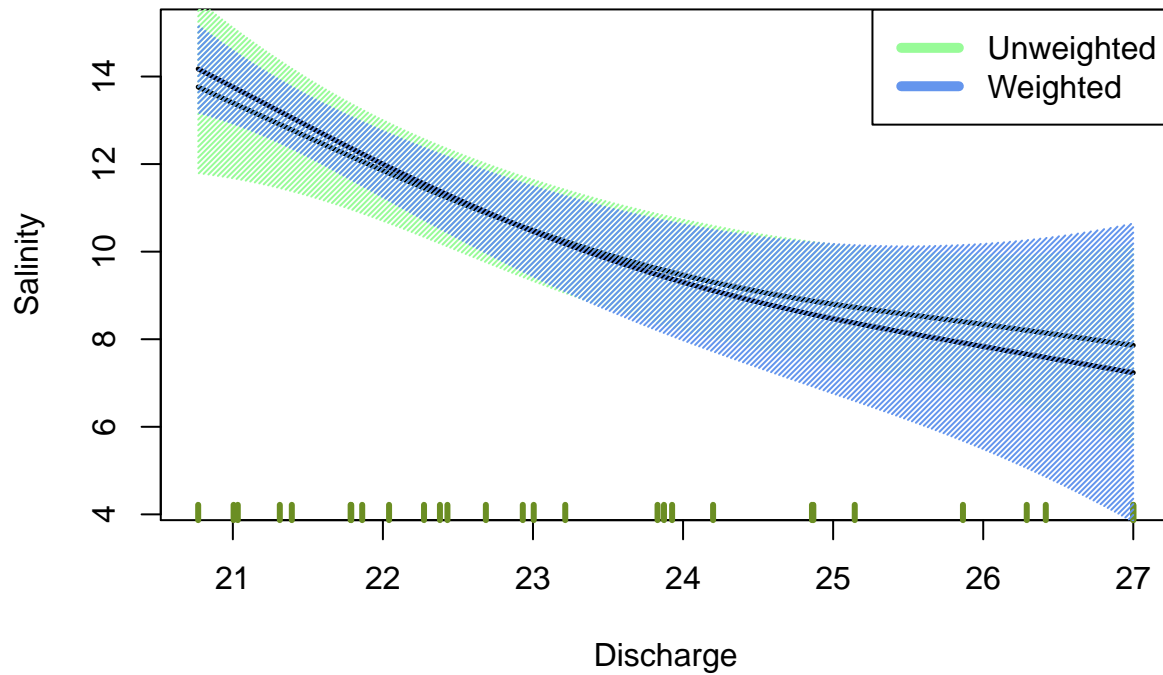
# Salinity by Discharge



I also found it helpful to see two side-by-side plots, due to the overlapping confidence bands in the above image.
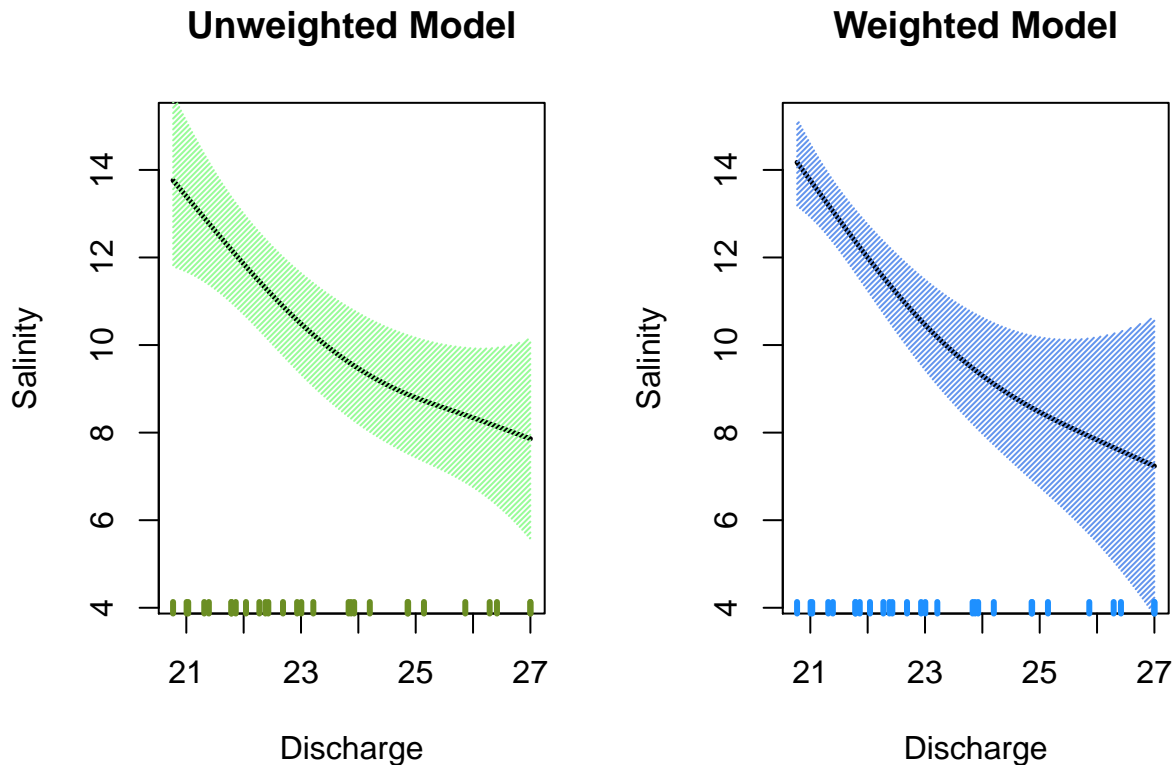
```r
par(mfrow=c(1,2))

plot(salinity$discharge, salinity$salinity, type="n"
     , main="Unweighted Model"
     , xlab="Discharge"
     , ylab="Salinity"
     )

# unweighted GAM first
gam_plotter(pred_obj=unweight_pred
            , lower_bnd=unweight_lower
            , upper_bnd=unweight_upper
            , shader="palegreen"
            )
rug(salinity$discharge, lwd=3, col="olivedrab")

plot(salinity$discharge, salinity$salinity, type="n"
     , main="Weighted Model"
     , xlab="Discharge"
     , ylab="Salinity"
     )

# weighted GAM second
gam_plotter(pred_obj=weight_pred
            , lower_bnd=weight_lower
            , upper_bnd=weight_upper
            , shader="cornflowerblue"
```

```
                    )
rug(salinity$discharge, lwd=3, col="dodgerblue")
```

**Unweighted Model**

**Weighted Model**



## Question 6

Are there any systematic differences in the two plots?

It appears that the weighted fit has generally narrower confidence bands in the regions of x where we have more observations. In the regions where we do not have as many observations (x>25), it appears that the confidence bands are wider.

## Part 2

## Question 7

For the Framingham data, fit a factor-by-curve GAM where CHD is regressed on LSBP, LCholest, and age (all splines), but allow an interaction between s(age) and smoker.

```
# first, read in the data
fram <- read.csv("/Users/panders2/Documents/schools/tamu/stat_689/homework/semiparametric-regression/mis
names(fram) <- tolower(names(fram))
str(fram)

## 'data.frame':    1615 obs. of  12 variables:
##  $ obs     : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ age     : int  56 38 54 42 47 43 58 43 43 36 ...
```

6

```
## $ sbp21   : int  110 118 140 122 136 138 108 130 116 150 ...
## $ sbp22   : int  120 126 120 116 134 120 98 120 110 146 ...
## $ sbp31   : int  136 120 125 115 130 118 112 112 108 122 ...
## $ sbp32   : int  134 114 124 110 144 110 105 120 102 115 ...
## $ smoker  : int  0 1 1 1 1 1 1 0 1 0 ...
## $ cholest2: int  276 211 284 225 292 192 194 319 205 247 ...
## $ cholest3: int  295 255 287 285 240 208 208 334 242 244 ...
## $ chd     : int  0 0 0 0 0 0 0 0 0 1 ...
## $ lsbp    : num  4.32 4.24 4.35 4.19 4.45 ...
## $ lcholest: num  5.65 5.45 5.65 5.54 5.58 ...
```

```r
fram_mod <- mgcv::gam(chd ~ s(lsbp, bs="cr", k=23) + s(lcholest, bs="cr", k=23)
                      + s(age, by=factor(smoker))
                      , family=binomial(link="logit")
                      , data=fram
                      )
summary(fram_mod)
```
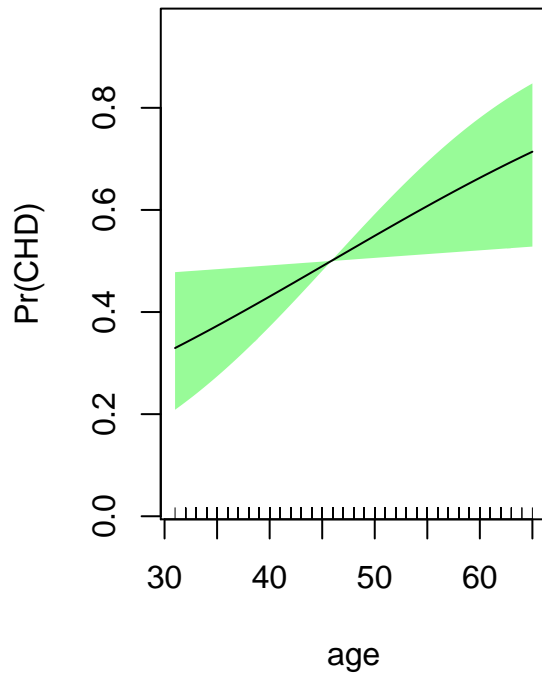
```
##
## Family: binomial
## Link function: logit
##
## Formula:
## chd ~ s(lsbp, bs = "cr", k = 23) + s(lcholest, bs = "cr", k = 23) +
##     s(age, by = factor(smoker))
##
## Parametric coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -2.7853     0.1222   -22.8   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##                          edf Ref.df Chi.sq  p-value
## s(lsbp)                1.623  2.044 15.579 0.000471 ***
## s(lcholest)            1.001  1.002 19.904 8.24e-06 ***
## s(age):factor(smoker)0 1.000  1.000  5.207 0.022507 *
## s(age):factor(smoker)1 2.937  3.658 20.239 0.000337 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.0449   Deviance explained = 9.38%
## UBRE = -0.48867  Scale est. = 1          n = 1615
```

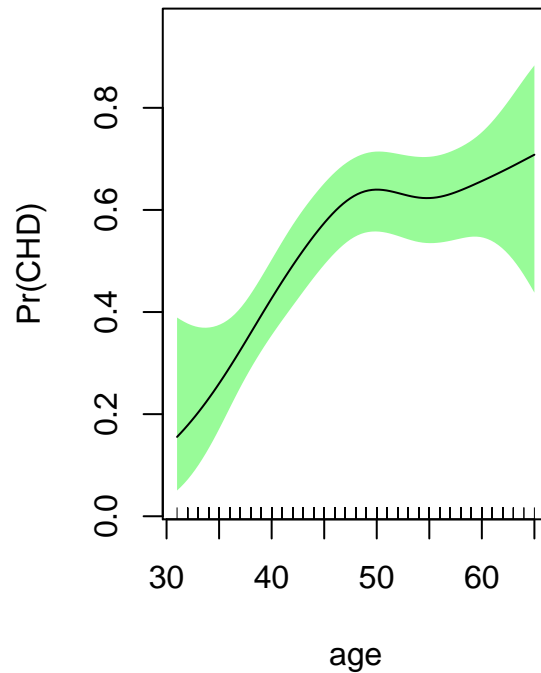Plot the fitted splines to age for smokers and non-smokers.

```r
par(mfrow=c(1,2))
plot.gam(fram_mod, shade = TRUE, shade.col="palegreen", select=3
         , trans=plogis
         , ylab="Pr(CHD)"
         , main="Smoker=0 by Age interaction")
plot.gam(fram_mod, shade = TRUE, shade.col="palegreen", select=4
         , trans=plogis
         , ylab="Pr(CHD)"
         , main="Smoker=1 by Age Interaction"
         )
```

## Smoker=0 by Age interaction

## Smoker=1 by Age Interaction



```r
par(mfrow=c(1,1))
```

# Question 8

In Question 7, do you think there is any statistical evidence that the risk of developing CHD over 8 years depends on whether or not you are a smoker?

For non-smokers, the partial dependence plot shows that the probability of CHD increases linearly, but for smokers, the max probability is achieved earlier.

```r
fram_mod$coefficients
```

```
##          (Intercept)            s(lsbp).1             s(lsbp).2
##        -2.785255e+00        -7.720861e-01         -6.351386e-01
##            s(lsbp).3            s(lsbp).4             s(lsbp).5
##        -2.861688e-01         3.491848e-02          2.515227e-01
##            s(lsbp).6            s(lsbp).7             s(lsbp).8
##         5.101281e-01         6.579107e-01          6.527723e-01
##            s(lsbp).9           s(lsbp).10            s(lsbp).11
##         6.645346e-01         6.884277e-01          6.569755e-01
##           s(lsbp).12           s(lsbp).13            s(lsbp).14
##         6.818382e-01         5.783954e-01          6.593683e-01
##           s(lsbp).15           s(lsbp).16            s(lsbp).17
##         6.570574e-01         6.970088e-01          6.894046e-01
##           s(lsbp).18           s(lsbp).19            s(lsbp).20
##         7.271983e-01         8.032995e-01          8.379079e-01
##           s(lsbp).21           s(lsbp).22         s(lcholest).1
##         9.808765e-01         1.227345e+00         -7.651977e-01
```

```
##           s(lcholest).2              s(lcholest).3              s(lcholest).4
##          -6.424044e-01              -4.225498e-01              -2.770511e-01
##           s(lcholest).5              s(lcholest).6              s(lcholest).7
##          -9.166235e-02               9.556440e-02               3.507158e-01
##           s(lcholest).8              s(lcholest).9             s(lcholest).10
##           4.818830e-01               5.400301e-01               6.172399e-01
##          s(lcholest).11             s(lcholest).12             s(lcholest).13
##           6.056508e-01               6.643179e-01               7.271522e-01
##          s(lcholest).14             s(lcholest).15             s(lcholest).16
##           7.306897e-01               7.470285e-01               7.086554e-01
##          s(lcholest).17             s(lcholest).18             s(lcholest).19
##           7.493413e-01               8.342322e-01               8.763964e-01
##          s(lcholest).20             s(lcholest).21             s(lcholest).22
##           8.990389e-01               1.193714e+00               2.158129e+00
## s(age):factor(smoker)0.1 s(age):factor(smoker)0.2 s(age):factor(smoker)0.3
##          -6.682757e-06               3.080497e-06               7.930601e-07
## s(age):factor(smoker)0.4 s(age):factor(smoker)0.5 s(age):factor(smoker)0.6
##           5.634168e-06               2.481209e-07               5.426805e-06
## s(age):factor(smoker)0.7 s(age):factor(smoker)0.8 s(age):factor(smoker)0.9
##           8.209319e-07               2.459328e-05               4.105063e-01
## s(age):factor(smoker)1.1 s(age):factor(smoker)1.2 s(age):factor(smoker)1.3
##           3.923161e-01              -2.988567e-01               4.272280e-03
## s(age):factor(smoker)1.4 s(age):factor(smoker)1.5 s(age):factor(smoker)1.6
##           1.362450e-01               3.309396e-02              -1.769286e-03
## s(age):factor(smoker)1.7 s(age):factor(smoker)1.8 s(age):factor(smoker)1.9
##          -1.645281e-02               3.516175e-01               9.216141e-01
```