

# Non-linear Regression and Classification

## Lab: Non-linear Modeling

In this lab, we analyse the `wage` data in order to illustrate the fact that many of the complex non-linear fitting procedures discussed can be easily implemented in R. We begin by loading the `ISLR2` library, which contains the data.

```
library(ISLR2)
attach(wage)
```

## Polynomial Regression

We first fit the model using the following command:

```
fit1 <- lm(wage ~ poly(age, 4), data = wage)
coef(summary(fit1))

##               Estimate Std. Error    t value    Pr(>|t|)
## (Intercept)    111.70361    9.7287669   11.328935  6.88880e-09
## poly(age, 4) 1  447.06785  39.9167851  11.289558  1.484604e-28
## poly(age, 4) 2 -478.21581  39.9167851 -11.983024  2.35951e-32
## poly(age, 4) 3  125.52169  39.9167851  3.144742  1.678622e-03
## poly(age, 4) 4 -77.81118  39.9167851 -1.951928  5.138365e-02
```

This syntax fits a linear model, using the `lm()` function, in order to predict `wage` using a fourth-degree polynomial in `age`: `poly(age, 4)`. The `poly()` command allows us to avoid having to write out a long formula with powers of `age`. The function returns a matrix whose columns are a basis of orthogonal polynomials, which essentially means that each column is a linear combination of the variables `age`, `age^2`, `age^3`, and `age^4`.

However, we can also use `poly()` to obtain `age`, `age^2`, `age^3` and `age^4` directly if we prefer. We can do this by using the `raw = TRUE` argument to the `poly()` function. Later we see that this does not affect the model in a meaningful way—though the choice of basis clearly affects the coefficient estimates, it does not affect the fitted values obtained.

```
fit2 <- lm(wage ~ poly(age, 4, raw = T), data = wage)
coef(summary(fit2))

##               Estimate Std. Error    t value    Pr(>|t|)
## (Intercept)    -1.841542e+02  6.094838e+01 -3.057172  8.6021802539
## poly(age, 4, raw = T) 1  2.124552e+01  5.888745e+00  3.689942  8.603123018
## poly(age, 4, raw = T) 2  6.826930e-01  2.861893e-01 -2.375743  8.602566446
## poly(age, 4, raw = T) 3  6.816688e-03  3.055931e-03  2.221469  8.6263977518
## poly(age, 4, raw = T) 4 -3.283383e-05  1.641559e-05 -1.951928  8.651386498
```

There are several other equivalent ways of fitting this model, which showcase the flexibility of the formula language in R. For example

```
fit2a <- lm(wage ~ age + I(age^2) + I(age^3) + I(age^4),
data = wage)
coef(fit2a)

## (Intercept)      age      I(age^2)      I(age^3)      I(age^4)
## -1.841542e+02  2.124552e+01 -5.838939e-01  6.816688e-03 -3.283383e-05
```

This simply creates the polynomial basis functions on the fly, taking care to protect terms like `age^2` via the function `I()` (the `^` symbol has a special meaning in formulas).

```
fit2b <- lm(wage ~ cbind(age, age^2, age^3, age^4),
data = wage)
```

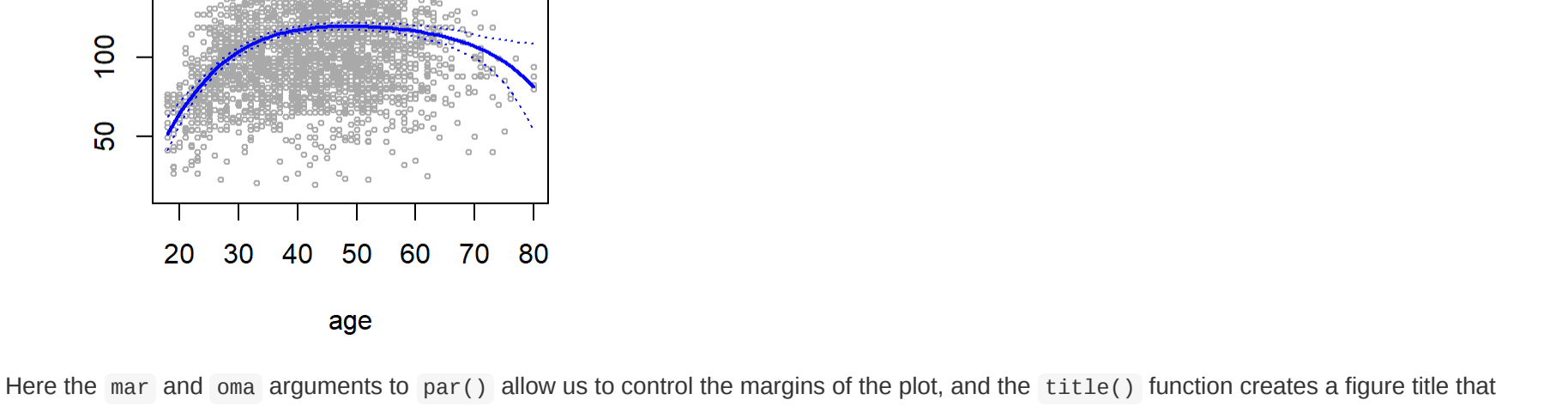
This does the same more compactly, using the `cbind()` function for building a matrix from a collection of vectors; any function call such as `cbind()` inside a formula also serves as a wrapper.

We now create a grid of values for `age` at which we want predictions, and then call the generic `predict()` function, specifying that we want standard errors as well.

```
age.lim <- range(age)
age.grid <- seq(from = age.lim[1], to = age.lim[2])
preds <- predict(fit1, newdata = list(age = age.grid),
se = TRUE)
se.bands <- cbind(preds$fit + 2 * preds$se.fit,
se$lower <- preds$fit - 2 * preds$se.fit)
```

Finally, we plot the data and add the fit from the degree-4 polynomial.

```
par(mfrow = c(1, 2), mar = c(4.5, 4.5, 1, 1),
oma = c(6, 4, 4, 8))
plot(age, wage, xlab = age.lim, cex = .5, col = "darkgrey")
title("Degree 4 Polynomial", outer = T)
lines(age.grid, preds$fit, lwd = 2, col = "blue")
matlines(age.grid, se.bands, lwd = 1, col = "blue", lty = 3)
```



Here the `mar` and `oma` arguments to `par()` allow us to control the margins of the plot, and the `title()` function creates a figure title that spans both subplots.

We mentioned earlier that whether or not an orthogonal set of basis functions is produced in the `poly()` function will not affect the model obtained in a meaningful way. What do we mean by this? The fitted values obtained in either case are identical.

```
preds2 <- predict(fit2, newdata = list(age = age.grid),
se = TRUE)
max(abs(preds$fit - preds2$fit))

## [1] 7.81597e-11
```

In performing a polynomial regression we must decide on the degree of the polynomial to use. One way to do this is by using hypothesis tests. We now fit models ranging from linear to a degree-5 polynomial and seek to determine the simplest model which is sufficient to explain the relationship between `wage` and `age`. We use the `anova()` function, which performs an ANOVA, using an *F*-test in order to test the null hypothesis that a model  $M_1$  is sufficient to explain the data against the alternative hypothesis that a more complex model  $M_2$  is required. In order to use the `anova()` function,  $M_1$  and  $M_2$  must be nested models: the predictors in  $M_1$  must be a subset of the predictors in  $M_2$ . In this case, we fit five different models and sequentially compare the simpler model to the more complex model.

```
fit.1 <- lm(wage ~ poly(age, 2), data = wage)
fit.2 <- lm(wage ~ poly(age, 3), data = wage)
fit.3 <- lm(wage ~ poly(age, 4), data = wage)
fit.4 <- lm(wage ~ poly(age, 5), data = wage)
fit.5 <- lm(wage ~ poly(age, 6), data = wage)
anova(fit.1, fit.2, fit.3, fit.4, fit.5)
```

```
## Analysis of Variance Table
##
## Model 1: wage ~ age
## Model 2: wage ~ poly(age, 2)
## Model 3: wage ~ poly(age, 3)
## Model 4: wage ~ poly(age, 4)
## Model 5: wage ~ poly(age, 5)
## Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1    2999 5622236    2 228786 143.5931 < 2.2e-16 ***
## 3    2996 4777674    1  15756   8.8888  0.001879 **
## 4    2996 4771086    1   6978   8.8898  0.001868 **
## 5    2996 4771822    1   1283   8.8968  0.366682
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The *p*-value comparing the linear `Model 1` to the quadratic `Model 2` is essentially zero ( $<10^{-16}$ ), indicating that a linear fit is not sufficient. Similarly the *p*-value comparing the quadratic `Model 2` to the cubic `Model 3` is very low (0.0017), so the quadratic fit is also insufficient. The *p*-value comparing the cubic and degree-4 polynomials, `Model 3` and `Model 4`, is approximately 5.5% while the degree-5 polynomial `Model 5` seems unnecessary because its *p*-value is 0.37. Hence, either a cubic or a quartic polynomial appear to provide a reasonable fit to the data, but lower- or higher-order models are not justified.

In this case, instead of using the `anova()` function, we could have obtained these *p*-values more succinctly by exploiting the fact that `poly()` creates orthogonal polynomials.

```
coef(summary(fit.5))

##               Estimate Std. Error    t value    Pr(>|t|)
## (Intercept)    111.70361    9.7287669   11.328935  6.88880e-09
## poly(age, 5) 1  447.06785  39.9167851  11.289558  1.481111e-28
## poly(age, 5) 2 -478.21581  39.9167851 -11.983024  2.357734e-32
## poly(age, 5) 3  125.52169  39.9167851  3.144632  1.679213e-03
## poly(age, 5) 4 -77.81118  39.9167851 -1.951874  5.186226e-02
## poly(age, 5) 5 -35.81389  39.9167851 -0.897248  3.69639e-01
```

Notice that the *p*-values are the same, and in fact the square of the *t*-statistics are equal to the *F*-statistics from the `anova()` function; for example:

```
(-11.983)^2

## [1] 143.5923
```

However, the ANOVA method works whether or not we used orthogonal polynomials; it also works when we have other terms in the model as well. For example, we can use `anova()` to compare these three models:

```
fit.1 <- lm(wage ~ education + age, data = wage)
fit.2 <- lm(wage ~ education + poly(age, 3), data = wage)
fit.3 <- lm(wage ~ education + poly(age, 3), data = wage)
anova(fit.1, fit.2, fit.3)
```

```
## Analysis of Variance Table
##
## Model 1: wage ~ education + age
## Model 2: wage ~ education + poly(age, 2)
## Model 3: wage ~ education + poly(age, 3)
## Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1    2994 3867892    2 142597 114.6969 <2e-16 ***
## 3    2992 3723896    1   5587   4.4936  0.0341 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

As an alternative to using hypothesis tests and ANOVA, we could choose the polynomial degree using cross-validation, as discussed in Chapter 5.

## Step functions

In order to fit a step function we use the `cut()` function.

```
table(cut(age, 4))

##
##      (17.9,33.5] (33.5,49] (49,64.5] (64.5,80.1]
##      750       1399       770       72

fit <- lm(wage ~ cut(age, 4), data = wage)
coef(summary(fit))

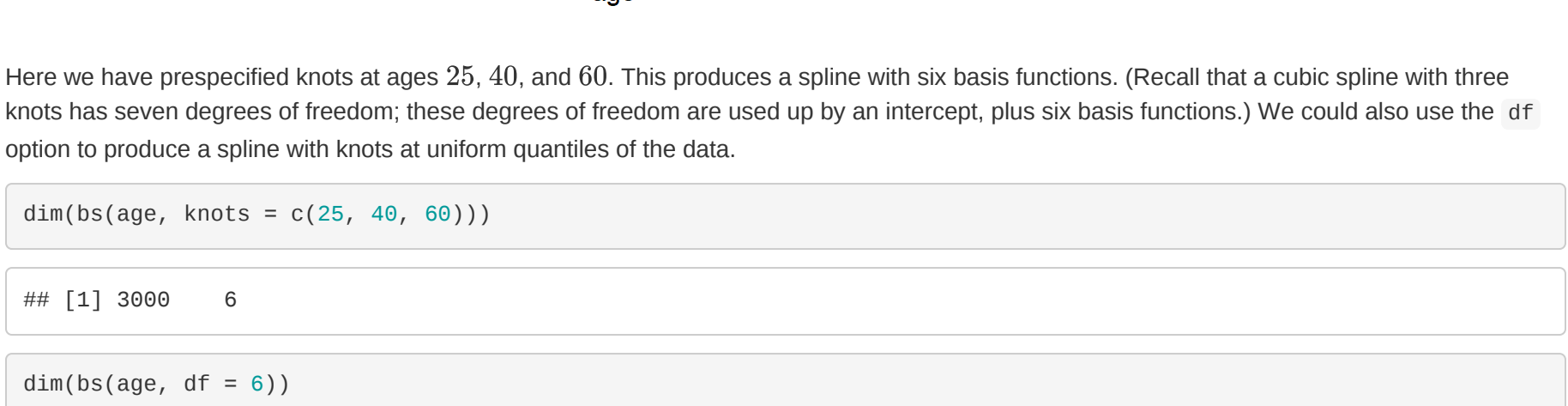
##               Estimate Std. Error    t value    Pr(>|t|)
## (Intercept)    64.153392   1.476569  43.79970  8.86880e-08
## cut(age, 4) (33.5,49]  24.853491   1.879431  13.148974  1.982315e-38
## cut(age, 4) (49,64.5]  23.864559   2.867558  8.144344  1.840758e-29
## cut(age, 4) (64.5,80.1]  7.846992   4.987424  1.53372  1.286388e-01
```

Here `cut()` automatically picked the cutpoints at 33.5, 49, and 64.5-years of age. We could also have specified our own cutpoints directly using the `breaks` option. The function `cut()` returns an ordered categorical variable: the `lm()` function then creates a set of dummy variables for use in the regression. The `age < 33.5` category is left out, so the intercept coefficient of \$64,160 can be interpreted as the average salary for those under 33.5-years of age, and the other coefficients can be interpreted as the average additional salary for those in the other age groups. We can produce predictions and plots just as we did in the case of the polynomial fit.

## Splines

In order to fit regression splines in R, we use the `splines` library. In Section 7.4, we saw that regression splines can be fit by constructing an appropriate matrix of basis functions. The `bs()` function generates the entire matrix of basis functions for splines with the specified set of knots. By default, cubic splines are produced. Fitting `wage` to `age` using a regression spline is simple:

```
library(splines)
fit <- lm(wage ~ bs(age, knots = c(25, 49, 69))), data = wage)
pred <- predict(fit, newdata = list(age = age.grid), se = T)
plot(age, wage, col = "gray")
lines(age.grid, preds$fit, lwd = 2)
lines(age.grid, preds$fit + 2 * preds$se, lty = "dashed", col = "red")
lines(age.grid, preds$fit - 2 * preds$se, lty = "dashed", col = "red")
```



Here we have prespecified knots at ages 25, 40, and 60. This produces a spline with six basis functions. (Recall that a cubic spline with three knots has seven degrees of freedom; these degrees of freedom are used up by an intercept, plus six basis functions.) We could also use the `df` option to produce a spline with knots at uniform quantiles of the data.

```
dim(bs(age, knots = c(25, 49, 60)))

## [1] 3009 6

dim(bs(age, df = 6))

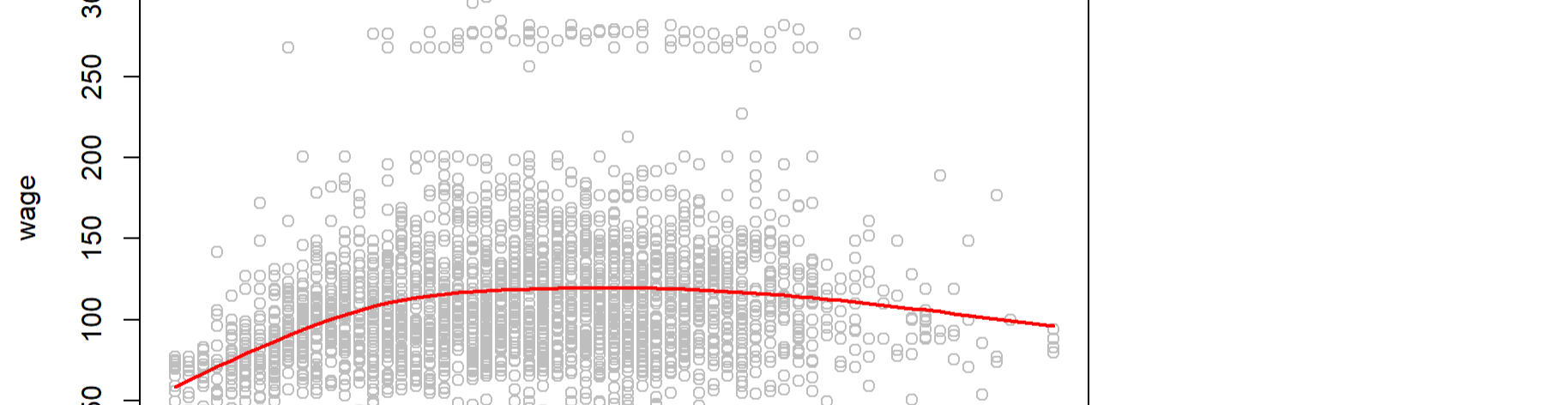
## [1] 3009 6

attr(,"knots")

##      25%      50%      75%
## 33.75 42.00 51.00
```

In this case R chooses knots at ages 33.8, 42.0, and 51.0, which correspond to the 25th, 50th, and 75th percentiles of `age`. The function `bs()` also has a `degree` argument, so we can fit splines of any degree, rather than the default degree of 3 (which yields a cubic spline).

```
fit2 <- lm(wage ~ ns(age, df = 4), data = wage)
pred2 <- predict(fit2, newdata = list(age = age.grid),
se = T)
plot(age, wage, col = "gray")
lines(age.grid, preds2$fit, col = "red", lwd = 2)
```



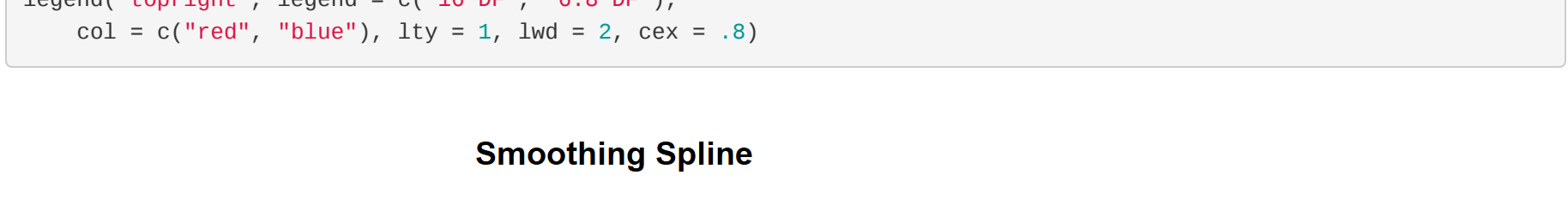
As with the `bs()` function, we could instead specify the knots directly using the `knots` option.

In order to fit a smoothing spline, we use the `smooth.spline()` function. Figure 7.8 was produced with the following code:

```
plot(age, wage, xlab = age.lim, cex = .5, col = "darkgrey")
title("Smoothing Spline")
fit <- smooth.spline(age, wage, df = 10)
fit2 <- smooth.spline(age, wage, cv = TRUE)
fit2$df

## [1] 6.794986

lines(fit, col = "red", lwd = 2)
lines(fit2, col = "blue", lwd = 2)
legend("topright", legend = c("df=10", "df=6.8 DF"),
col = c("red", "blue"), lty = 1, lwd = 2, cex = .8)
```



Notice that in the first call to `smooth.spline()`, we specified `df = 16`. The function then determines which value of  $\lambda$  leads to 16 degrees of freedom, in the second call to `smooth.spline()`, we select the smoothness level by cross-validation; this results in a value of  $\lambda$  that yields 6.8 degrees of freedom.

## GAMs

We now fit a GAM to predict `wage` using natural spline functions of `year` and `age`, treating `education` as a qualitative predictor, as in (7.16). Since this is just a big linear regression model using an appropriate choice of basis functions, we can simply do this using the `lm()` function.

```
gam1 <- lm(wage ~ s(year, 4) + ns(age, 5) + education,
data = wage)
```

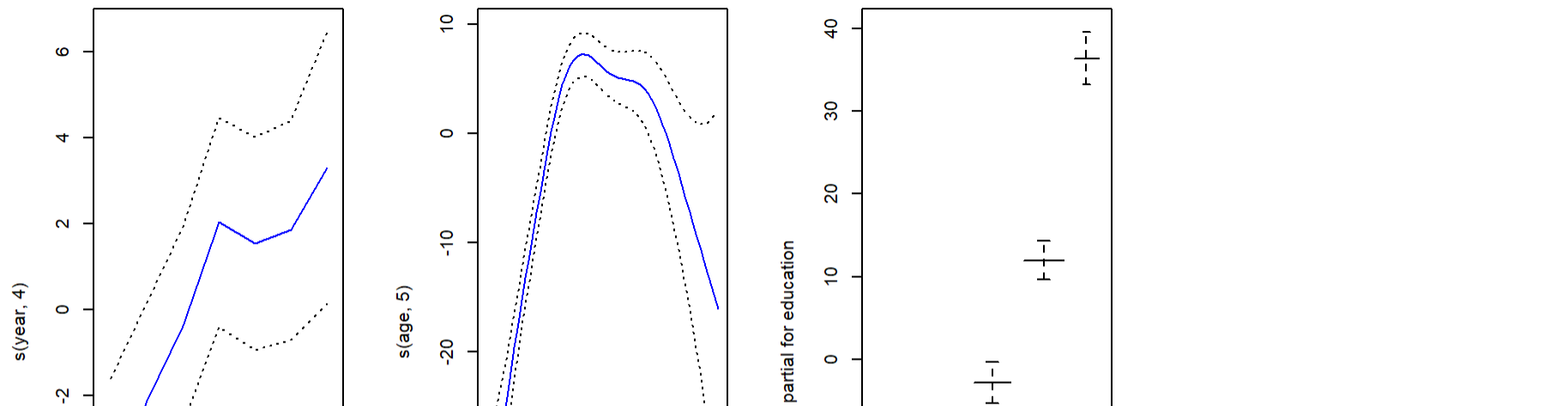
We now fit the model (7.16) using smoothing splines rather than natural splines. In order to fit more general sorts of GAMs, using smoothing splines or other components that cannot be expressed in terms of basis functions and then fit using least squares regression, we will need to use the `gam` library in R.

The `s()` function, which is part of the `gam` library, is used to indicate that we would like to use a smoothing spline. We specify that the function of `year` should have 4 degrees of freedom, and that the function of `age` will have 5 degrees of freedom. Since `education` is qualitative, we leave it as is, and it is converted into four dummy variables. We use the `gam()` function in order to fit a GAM using these components. All of the terms in (7.16) are fit simultaneously, taking each other into account to explain the response.

```
library(gam)
gam.m3 <- gam(wage ~ s(year, 4) + s(age, 5) + education,
data = wage)
```

In order to produce Figure 7.12, we simply call the `plot()` function:

```
plot(gam.m3, se = TRUE, col = "blue")
```



The generic `plot()` function recognizes that `gam.m3` is an object of class `gam`, and invokes the appropriate `plot.gam()` method. Conveniently, even though `gam1` is not of class `gam` but rather of class `lm`, we can (just use `plot.gam()` on it. Figure 7.11 was produced using the following expression:

```
plot(gam1, gam1, se = TRUE, col = "red")
```



Notice here we had to use `plot.gam()` rather than the generic `plot()` function.

In these plots, the function of `year` looks rather linear. We can perform a series of ANOVA tests in order to determine which of these three models is best: a GAM that excludes `year` ( $M_1$ ), a GAM that uses a linear function of `year` ( $M_2$ ), or a GAM that uses a spline function of `year` ( $M_3$ ).

```
gam.m1 <- gam(wage ~ s(age, 5) + education, data = wage)
gam.m2 <- gam(wage ~ year + s(age, 5) + education,
data = wage)
anova(gam.m1, gam.m2, gam.m3, test = "F")
```

```
## Analysis of Variance Table
##
## Model 1: wage ~ s(age, 5) + education
## Model 2: wage ~ year + s(age, 5) + education
## Model 3: wage ~ s(year, 4) + s(age, 5) + education
## Resid. Df Resid. Dev Df Sum of Sq    F    Pr(>F)
## 1    2988  368977.3    1 17889.2 14.4771 0.988147 ***
## 2    2989  368984.2    1  1332.9 10.9339 356.081 < 2.2e-16 ***
## 3    2988  368977.3    4  4071.1 1.6982 0.3485661
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We find that there is compelling evidence that a GAM with a linear function of `year` is better than a GAM that does not include `year` at all. However, there is no evidence that a non-linear function of `year` is needed (*p*-value=0.349). In other words, based on the results of this ANOVA,  $M_2$  is preferred.

The `summary()` function produces a summary of the gam fit.

```
summary(gam.m3)

##
## Call: gam(formula = wage ~ s(year, 4) + s(age, 5) + education, data = wage)
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -119.43   -10.70    -5.33   14.17   233.48
## (Dispersion Parameter for gaussian family taken to be 1235.69)
##
## Null Deviance: 5228866 on 2989 degrees of freedom
## Residual Deviance: 3689778 on 2986 degrees of freedom
## AIC: 29897.75
##
## Number of Local Scoring Iterations: NA
##
## Aova for Parametric Effects
##
## Df Sum Sq Mean Sq F value    Pr(>F)
## s(year, 4)    1 27362 27362 21.982 2.877e-06 ***
## s(age, 5)    1 135329 135329 356.081 < 2.2e-16 ***
## education    4 3869726 2917432 236.421 < 2.2e-16 ***
## Residuals    2988 3689778 1236
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The ANOVA for Parametric Effects *p*-values clearly demonstrate that `year`, `age`, and `education` are all highly statistically significant, even when assuming a linear relationship. Alternatively, the ANOVA for Nonparametric Effects *p*-values for `year` and `age` correspond to a null hypothesis of a linear relationship versus the alternative of a non-linear relationship. The large *p*-value for `year` reinforces our conclusion from the ANOVA test that a linear function is adequate for this term. However, there is very clear evidence that a non-linear term is required for `age`.

We can make predictions using the `predict()` method for the class `gam`. Here we make predictions on the training set.

```
preds <- predict(gam.m2, newdata = wage)
```