

# hw3

*Xingwen Wei, Xin Hu, Liding Li*

*March 6, 2021*

## Question 1

The true model is a polynomial of degree 3.

```
m <- matrix(c('high', 'low', 'low', 'low', 'low', 'high'), ncol=2, byrow=FALSE)
colnames(m) <- c('Bias', 'Variance')
rownames(m) <- c('Linear regression', 'Polynomial regression with degree 3', 'Polynomial regression with degree 10')
as.table(m)
```

```
##                                Bias Variance
## Linear regression                high low
## Polynomial regression with degree 3  low low
## Polynomial regression with degree 10 low  high
```

## Question 2

a

As  $\lambda \rightarrow \infty$ ,  $\hat{g}_1$  will have all  $g^{(3)}(x) = 0$  and  $\hat{g}_2$  will have all  $g^{(4)}(x) = 0$ . So this is similar to constraining  $\hat{g}_1$  to have degree less than 3 and  $\hat{g}_2$  less than 4. Thus,  $\hat{g}_2$  will always have smaller or equal training error than  $\hat{g}_1$ .

b

On one hand, if the true curve has degree higher than or equal to 3,  $\hat{g}_1$  will not be able to capture it at all, while  $\hat{g}_2$  can capture it. So  $\hat{g}_2$  will have the smaller test error in this case. On the other hand, if the true curve has degree smaller than 3,  $\hat{g}_2$  may pick some noise up as signal and overfits the training data, while  $\hat{g}_1$  will not. So  $\hat{g}_1$  will have the smaller test error in this case.

c

For  $\lambda = 0$ ,  $\hat{g}_1 = \hat{g}_2$ . So they will have the same training and test error.

## Question 3

a

Exploratory Data Analysis

```
library(ggplot2)
library(tidyr)
library(GGally)

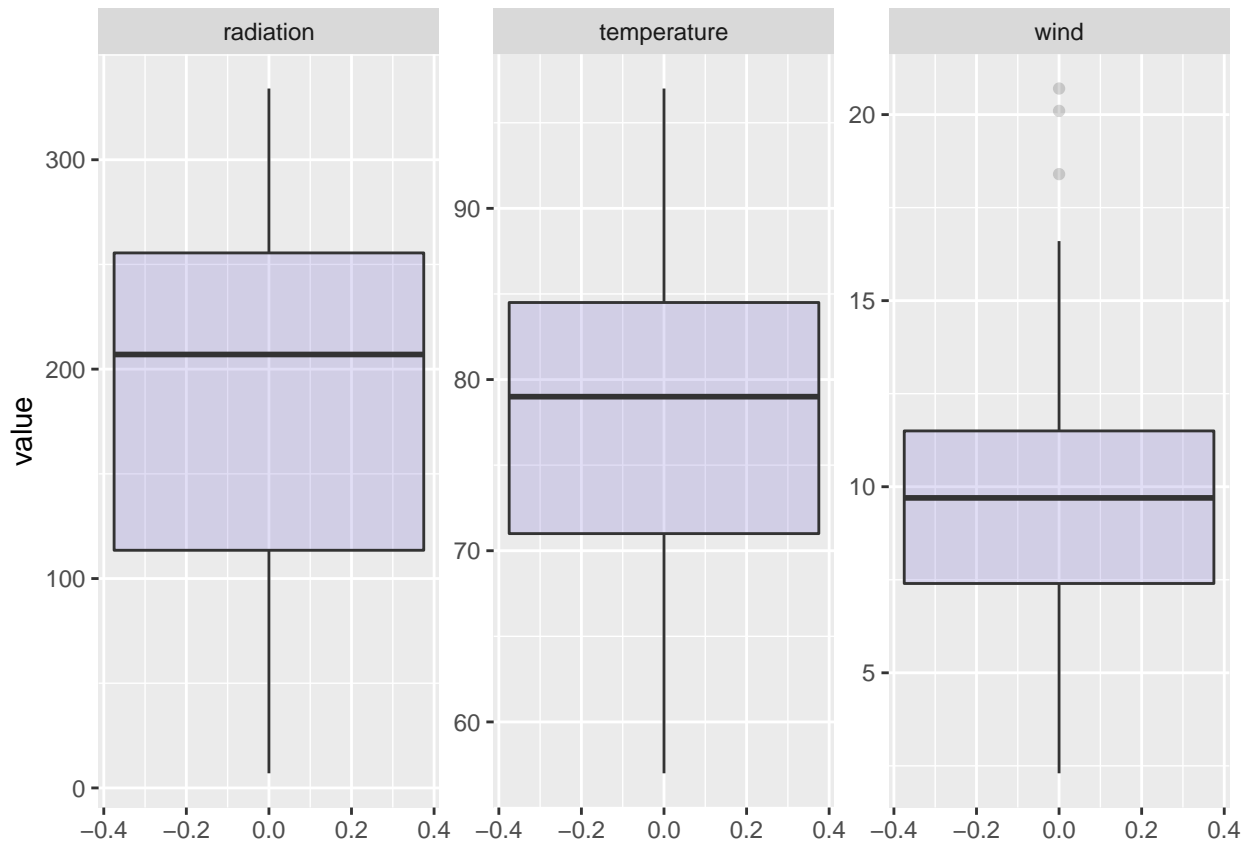
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2

ozone <- read.table("C:/Users/xingw/Desktop/503/stats503/hw3/ozone_data.txt", header=1)
summary(ozone)
```

```
##      ozone      radiation      temperature      wind
## Min.   :  1.0   Min.     :  7.0   Min.     :57.00   Min.     :  2.300
```

```
## 1st Qu.: 18.0    1st Qu.:113.5    1st Qu.:71.00    1st Qu.: 7.400
## Median : 31.0    Median :207.0    Median :79.00    Median : 9.700
## Mean   : 42.1    Mean   :184.8    Mean   :77.79    Mean   : 9.939
## 3rd Qu.: 62.0    3rd Qu.:255.5    3rd Qu.:84.50    3rd Qu.:11.500
## Max.   :168.0    Max.   :334.0    Max.   :97.00    Max.   :20.700
```

```
ozone_long <- gather(ozone, predictor, value, 2:4, factor_key = TRUE)
ggplot(ozone_long, aes(y=value)) + geom_boxplot(fill='slateblue', alpha=0.2) + facet_wrap(vars(predictor))
```



According to the model summary, the linear model we choose is

$$\text{ozone}^{\frac{1}{3}} = 0.001 \times \text{radiation} + 0.058 \times \text{temperature} - 0.066 \times \text{wind} - 0.852$$

```
set.seed(123)
ozone <- read.table("C:/Users/xingw/Desktop/503/stats503/hw3/ozone_data.txt", header=1)
train <- sample(nrow(ozone), floor(nrow(ozone)*0.7))
ozone$cbr <- ozone$ozone^(1/3)
lmod <- lm(cbr~radiation+temperature+wind, data=ozone[train, ])
summary(lmod)
```

```
##
## Call:
## lm(formula = cbr ~ radiation + temperature + wind, data = ozone[train,
##      ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -0.94503 -0.40230 -0.00071 0.27566 1.50475
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.8521380  0.6449384  -1.321 0.190538
## radiation    0.0016477  0.0006398   2.575 0.012037 *
## temperature  0.0579790  0.0071750   8.081 9.93e-12 ***
## wind         -0.0656469  0.0180658  -3.634 0.000516 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4999 on 73 degrees of freedom
## Multiple R-squared:  0.7111, Adjusted R-squared:  0.6992
## F-statistic: 59.9 on 3 and 73 DF,  p-value: < 2.2e-16
```

b

We use LOOCV to find to optimal number of knots. According to the plot below, we find that we get the best result when the number of knots is 2. Based on the fitted splines plot of each variables, we find that temperature is linear where radiation and wind are not.

```
library(gam)
```

```
## Warning: package 'gam' was built under R version 3.6.3
```

```
## Loading required package: splines
```

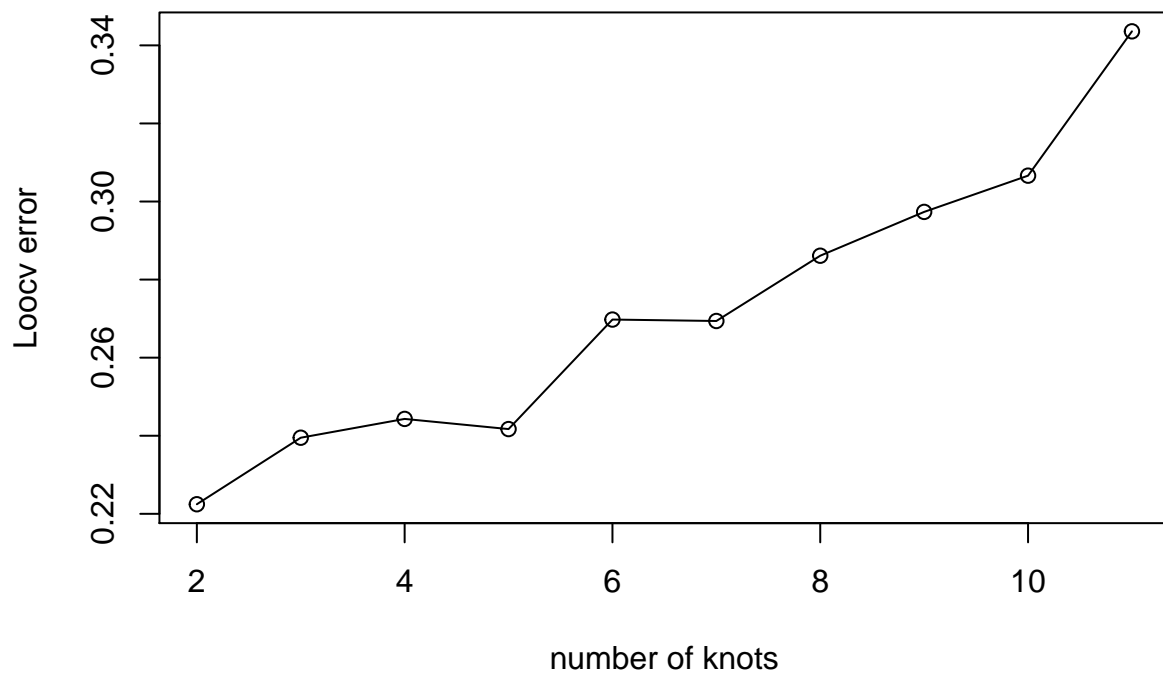
```
## Loading required package: foreach
```

```
## Warning: package 'foreach' was built under R version 3.6.3
```

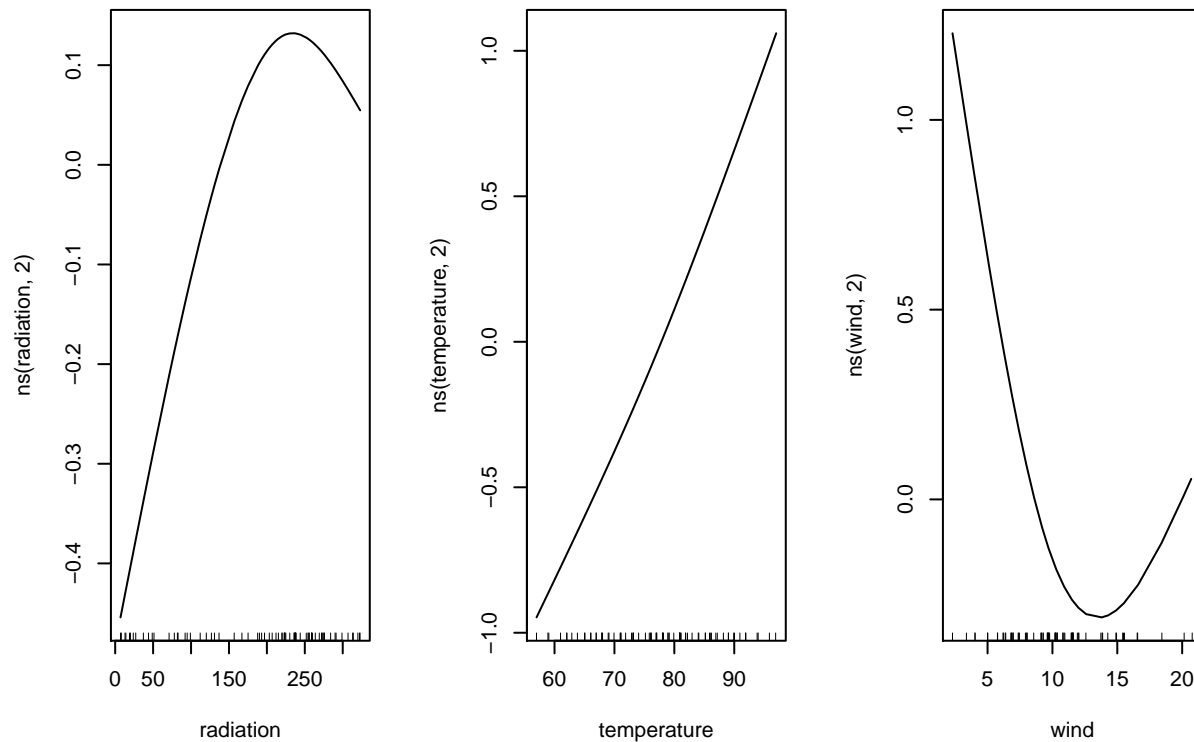
```
## Loaded gam 1.20
```

```
#number_of_knots <- 2:10
train_loo <- ozone[train, ]
errors <- matrix(NA, nrow=nrow(train_loo), ncol=10)

for(i in 1:nrow(train_loo)){
  train1 <- train_loo[-i, ]
  test1 <- train_loo[i, ]
  for(k in 1:10){
    #q <- seq(from=0, to=1, length.out=knot+2)
    #q <- q[2:(length(q)-1)]
    gam_mod <- gam(cbr~ns(radiation, (k+1))+ns(temperature, (k+1))+ns(wind, (k+1)), data=train1)
    pred <- predict(gam_mod, test1)
    errors[i, k] <- (test1$cbr-pred)^2
  }
}
plot(x=2:11, y=apply(errors, 2, mean), 'o', xlab='number of knots', ylab='Loocv error')
```



```
gam_mod <- gam(cbr~ns(radiation, 2)+ns(temperature, 2)+ns(wind, 2), data=ozone[train, ])  
par(mfrow=c(1, 3))  
plot.Gam(gam_mod)
```

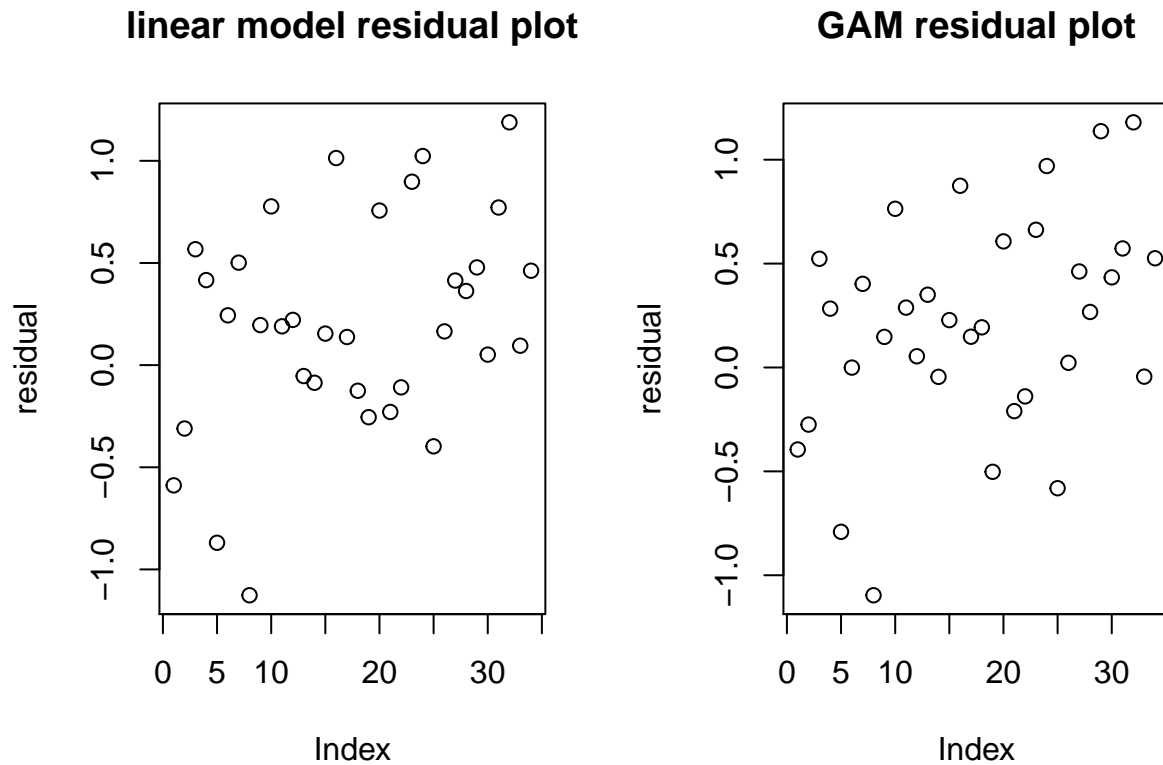


c

The mean test error for linear model is 0.312. On the other hand, we get a test error of 0.306 with GAM. According to both residual plots, the error terms are roughly normally distributed around 0 with no apparent pattern. We believe the similar result from both methods indicates the additional knots may not be too helpful. We suspect that the non-linearity we observed in fitted spline plots above is a result of small sample size.

```
par(mfrow=c(1, 2))
pred <- predict(lmod, ozone[-train, ])
lm_mse <- sum((pred-ozone[-train, 'cbr'])^2)/nrow(ozone[-train, ])
plot(pred-ozone[-train, 'cbr'], ylab='residual', main='linear model residual plot')

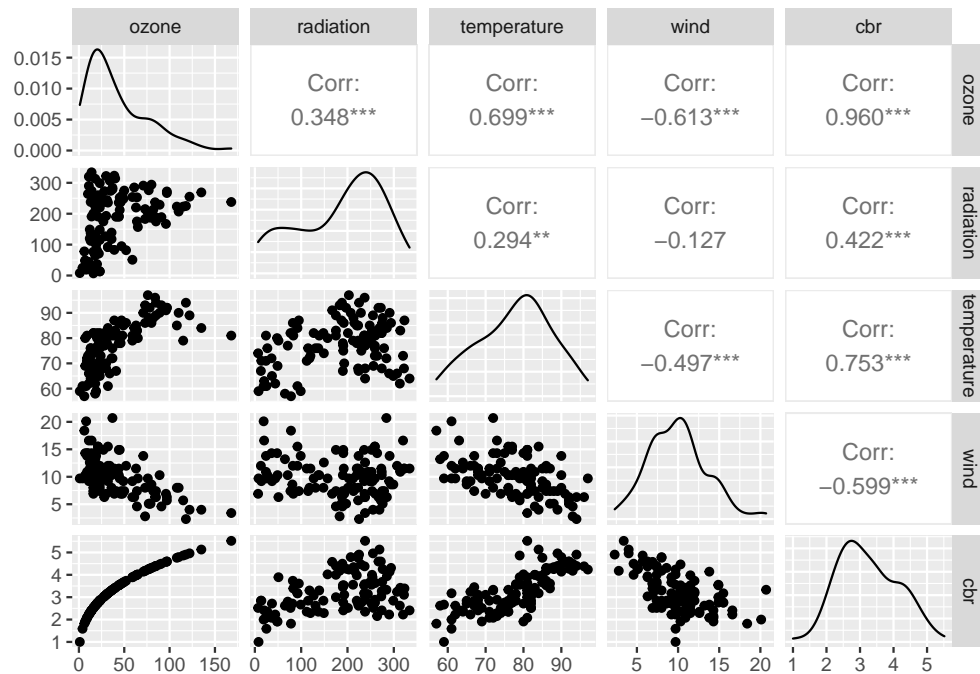
gam_pred <- predict(gam_mod, ozone[-train, ])
gam_mse <- sum((gam_pred-ozone[-train, 'cbr'])^2)/nrow(ozone[-train, ])
plot(gam_pred-ozone[-train, 'cbr'], ylab='residual', main='GAM residual plot')
```



**d**

According to the pairwise scatterplot, we find there is a rather strong linear relationship between temperature and wind and the cubic root of ozone, as corroborated by the correlation coefficients 0.75 and -0.6 respectively. One can argue that the radiation does not have a strong linear relationship with the response variable, with correlation coefficient of 0.42. However, it is clear that lower values of radiations are associated with lower values of cubic root of ozone and higher values of radiations are associated with higher values of cubic root of ozone. Therefore, we would not apply GAM in this dataset with about 100 observations without further strong evidence for nonlinearity.

```
ggpairs(ozone)
```



```
plot(ozone$radiation, ozone$cbr, xlab='radiation', ylab=expression(ozone^(1/3)))
ts <- lm(ozone$cbr~ozone$radiation)
pred1 <- predict(ts, ozone['radiation'])
lines(x=ozone$radiation, y=pred1, col='red')
```

