

Assignment 2, STK-INF 4300

He Gu

heg@mail.uio.no

Problem 1

Task 1

At the very first, we could check the distribution of each modality for all categorical variables.

CODE:

```
1. data <- read.table("C:/Users/he322/Desktop/Oblig2STK/ozone_496obs_25vars.txt", header = TRUE)
2. header <- names(data)
3. variables <- header[-25]
4. cata_vari <- c("ADHEU", "SEX", "HOCHOZON", "AMATOP", "AVATOP", "ADEKZ", "ARAUCH",
5.               "FSNIGHT", "FMILB", "FTIER", "FPOLL", "FLTOTMED", "FSPT", "FSATEM",
6.               "FSAUGE", "FSPFEI", "FSHLAUF")
7. num_vari <- c("ALTER", "AGEBGEW", "FLGROSS", "FNOH24", "FO3H24", "FTEH24", "FLGEW")
8. #1
9. #get and split and standardize the data
10. ##check the distribution for catagorical variables
11. print(c("==1", "==0"))
12. for (v in cata_vari){
13.   a0 <- nrow(data[data[v] == "1",])
14.   a1 <- nrow(data[data[v] == "0",])
15.   print(c(a0,a1))
16. }
```

RESULT:

```

1. [1] "==1" "==0"
2. [1] 37 459
3. [1] 251 245
4. [1] 315 181
5. [1] 132 364
6. [1] 104 392
7. [1] 97 399
8. [1] 150 346
9. [1] 55 441
10. [1] 79 417
11. [1] 56 440
12. [1] 120 376
13. [1] 47 445
14. [1] 154 342
15. [1] 26 470
16. [1] 68 428
17. [1] 26 470
18. [1] 33 463

```

Obviously, we see that there are clearly some unbalance, but since we are having 496 observations in total, it is nearly impossible that our training set contains only one modality even if we just randomly split the data. And we could test it.

We then simply split the data randomly and check the distribution of each modality.

CODE:

```

1. ##split the data randomly
2. n_random = sample(1:496,496)
3. train_set <- data[c(sort(n_random[1:248])),]
4. test_set <- data[c(sort(n_random[249:496])),]
5.
6. ##check the distribution for catagorical variables of both test and train sets
7. for (v in cata_vari){
8.   a0 <- nrow(train_set[train_set[v] == "1",])
9.   a1 <- nrow(train_set[train_set[v] == "0",])
10.  print(c(a0,a1))
11. }

```

RESULT:

```

1. [1] "==1" "==0"
2. [1] 21 227
3. [1] 126 122
4. [1] 157 91

```

```
5. [1] 57 191
6. [1] 46 202
7. [1] 47 201
8. [1] 73 175
9. [1] 29 219
10. [1] 41 207
11. [1] 29 219
12. [1] 64 184
13. [1] 29 218
14. [1] 83 165
15. [1] 10 238
16. [1] 38 210
17. [1] 12 236
18. [1] 15 233
```

We see that the training set does contain all modalities.

And we shall then standardize the data. Notice that there are categorical variables, but according to our tasks' description, we do not need to interpret our model, which means that these variables could also be standardized, although it might be pretty meaningless. But notice that we need to find the "most relevant" variable in further tasks, and it means that we are supposed to standardize also the categorical variables, if not, the relative "importance" of these categorical variables could be higher than they really are.

Thus, we are standardizing categorical variables by the same mechanism as standardizing the other variables, i.e. *Z – score normalization*

Since we are assuming that we have no idea about the test set, we need to standardize both training and test set by using the status of training set, i.e. both by

$$new\ data = \frac{data - mean(training\ set)}{std(training\ set)}$$

CODE:

```
1. ##standardize the training data and test set BOTH by mean(train) and std(train)
2. for (v in variables){
3.   mean_train <- mean(train_set[,v])
4.   sd_train <- sd(train_set[,v])
5.   train_set[v] <- (train_set[v] - mean_train)/ sd_train
6.   test_set[v] <- (test_set[v] - mean_train)/ sd_train
7. }
```

Task 2

CODE:

```
1. #2
2. #fit the data with gaussian
3. model <- lm(train_set$FFVC ~ ., data = train_set)
4. summary(model)
```

RESULT:

```
1. Coefficients:
2.              Estimate Std. Error t value Pr(>|t|)
3. (Intercept)  2.2727419  0.0128286 177.162 < 2e-16 ***
4. ALTER        0.0506641  0.0157764   3.211  0.00152 **
5. ADHEU       -0.0243432  0.0159593  -1.525  0.12859
6. SEX         -0.1029461  0.0136156 -7.561 1.03e-12 ***
7. HOCHOZON    -0.0091874  0.0180470  -0.509  0.61120
8. AMATOP       0.0007830  0.0138945   0.056  0.95511
9. AVATOP       0.0008852  0.0135605   0.065  0.94801
10. ADEKZ      -0.0124366  0.0140509  -0.885  0.37705
11. ARAUCH     -0.0029869  0.0139206  -0.215  0.83030
12. AGEBGEW     0.0046095  0.0139171   0.331  0.74080
13. FSNIGHT    -0.0017677  0.0149890  -0.118  0.90623
14. FLGROSS     0.1515287  0.0222518   6.810 8.96e-11 ***
15. FMILB      -0.0108766  0.0191260  -0.569  0.57015
16. FNOH24     -0.0176679  0.0179608  -0.984  0.32633
17. FTIER      -0.0066234  0.0159679  -0.415  0.67869
18. FPOLL      -0.0484411  0.0279185  -1.735  0.08411 .
19. FLTOTMED   -0.0126885  0.0138921  -0.913  0.36204
20. FO3H24      0.0431497  0.0286550   1.506  0.13352
21. FSPT        0.0260766  0.0314917   0.828  0.40853
22. FTEH24     -0.0428080  0.0268900  -1.592  0.11281
23. FSATEM      0.0125539  0.0164625   0.763  0.44652
24. FSAUGE     -0.0028839  0.0153242  -0.188  0.85090
25. FLGEW       0.0820594  0.0205223   3.999 8.67e-05 ***
26. FSPFEI      0.0453508  0.0162609   2.789  0.00575 **
27. FSHLAUF     0.0012133  0.0143145   0.085  0.93253
28. ---
29. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
30.
31. Residual standard error: 0.202 on 223 degrees of freedom
32. Multiple R-squared:  0.6853,    Adjusted R-squared:  0.6514
33. F-statistic: 20.23 on 24 and 223 DF,  p-value: < 2.2e-16
```

COMMENTS:

Apparently, we see that **FLGROSS** has the strongest association with the forced vital capacity in such case, since the absolute value of its estimate is higher than others.

As for the **standard error**, we see that many variables have a relatively small standard error, which means that their estimates could be stable. But there also exists some variables which have a relatively large standard error, and it means that the estimates of these variables could be quite **unstable**, e.g. **ADHEU**, **HOCHOZON**, and etc.

Besides, according to the *p*-values of our variables, some of them are significantly non-zero, i.e. **ALTER**, **SEX**, **FLGROSS**, **FLGEW** and **FSPFEI**, while others are not.

Task 3

At first, we are using **AIC** as the stopping criterion.

CODE:

```
1. library(MASS)
2. null_model <- lm(train_set$FFVC ~ 1, data = train_set)
3.
4. backward_AIC <- stepAIC(model, direction = "backward", k = 2, trace = FALSE)
5. forward_AIC <- stepAIC(null_model, direction = "forward", k = 2, trace = FALSE,
6.                        scope = list(lower = null_model, upper = model))
7. backward_AIC
8. forward_AIC
```

RESULT:

```
1. > (backward_AIC)
2.
3. Call:
4. lm(formula = FFVC ~ ALTER + SEX + HOCHOZON + ARAUCH + AGEBGEW +
5.     FLGROSS + FLGEW + FSPFEI, data = train_set)
6.
7. Coefficients:
8. (Intercept) ALTER SEX HOCHOZON ARAUCH AGEBGEW FLGROSS FLGEW FSPFEI
9. 2.29084 0.02265 -0.09708 -0.02270 0.02954 0.02528 0.19083 0.06379 0.03864
10.
11. > (forward_AIC)
12.
13. Call:
14. lm(formula = train_set$FFVC ~ FLGROSS + SEX + FLGEW + FSPFEI +
15.     ARAUCH + HOCHOZON + AGEBGEW + ALTER, data = train_set)
```

```

16.
17. Coefficients:
18. (Intercept)  FLGROSS    SEX    FLGEW    FSPFEI    ARAUCH    HOCHOZON    AGEBGEW    ALTER
19. 2.29084    0.19083   -0.09708    0.06379    0.03864    0.02954   -0.02270    0.02528    0.02265

```

COMMENTS:

Obviously, we see that the models selected by backward and forward AIC selection are the same, but in fact, they are not always the same.

Sometimes, the model selected by backward AIC may contain more variables than forward AIC. This is pretty reasonable, since the selection made by AIC may converge at a local minimum, which means that backward and forward selections might converge earlier than they supposed to be, and thus, the backward AIC contains more variables than backward AIC selection.

But in such case, it seems that our selection performs pretty well, since both backward AIC and forward AIC converge at the same model.

We then change the stopping criterion into **BIC**, i.e. change $k = 2$ into $k = \log(n) = \log(24)$ in such case, and run the similar code again.

CODE:

```

1. backward_BIC <- stepAIC(model, direction = "backward", k = log(24), trace = FALSE)
2. forward_BIC <- stepAIC(null_model, direction = "forward", k = log(24), trace = FALSE,
   scope = list(lower = null_model, upper = model))
3. backward_BIC
4. forward_BIC

```

RESULT:

```

1. > (backward_BIC)
2.
3. Call:
4. lm(formula = FFVC ~ SEX + FLGROSS + FPOLL + FLTOTMED + F03H24 +
5.     FTEH24 + FSATEM + FLGEW, data = train_set)
6.
7. Coefficients:
8. (Intercept)  SEX    FLGROSS    FPOLL    FLTOTMED    F03H24    FTEH24    FSATEM    FLGEW
9.  2.28480   -0.09755   0.17628   -0.03219  -0.02928   0.06308  -0.07099   0.04508   0.10742
10.
11. > (forward_BIC)
12.
13. Call:
14. lm(formula = train_set$FFVC ~ FLGROSS + SEX + FLGEW + FSATEM +
15.     FLTOTMED + FPOLL, data = train_set)
16.

```

17. Coefficients:

18. (Intercept)	FLGROSS	SEX	FLGEW	FSATEM	FLTOTMED	FPOLL
19. 2.28480	0.17377	-0.09719	0.10694	0.04752	-0.03165	-0.03201

COMMENTS:

Following the same logic we used to analyze AIC selection, we see that, this time, the backward selection does contain more variables than the forward selection.

We know that the forward selection starts from a null model while the backward one starts from a full model, and in such case, they converge at different local minimums, and this is why the selected models are different and backward selection contains more variables.

OVERALL COMMENTS:

In such case, we see that the backward selection contains more variables than the forward selection, and the AIC selection contains more variables than the BIC selection. This is exactly the same as we expected.

According to the formulas of AIC and BIC, we see that BIC penalizes the amount of variables, especially when the amount is relatively small, while AIC does not penalize the amount.

In addition, we know that AIC tends to find a model which performs better in prediction, since AIC assumes that there is no correct model among the candidates, and it shall only find the “closest” model to the true model. But meanwhile, BIC tends to find a sparser model, since it shall find the model which is “most likely” to be true.

In fact, it is hard to say which stopping criterion and which kind of procedure is better, but we could just simply test their performance with using the test set, and it shall be fulfilled in Task 7.

REPORT FOR THESE 4 MODELS:

CODE:

```
1. summary(backward_AIC)
2. summary(forward_AIC)
3. summary(backward_BIC)
4. summary(forward_BIC)
```

RESULT:

```
1. > summary(backward_AIC)
2.
3. Call:
4. lm(formula = FFVC ~ SEX + AMATOP + FLGROSS + FPOLL + FLTOTMED +
5.     F03H24 + FTEH24 + FSATEM + FLGEW + FSPFEI, data = train_set)
6.
7. Residuals:
```

```

8.      Min      1Q   Median      3Q      Max
9. -0.47336 -0.13900 -0.00203  0.13173  0.55450
10.
11. Coefficients:
12.      Estimate Std. Error t value Pr(>|t|)
13. (Intercept)  2.28480    0.01266 180.534 < 2e-16 ***
14. SEX         -0.09756    0.01321  -7.386 2.51e-12 ***
15. AMATOP       -0.02142    0.01324  -1.618  0.10703
16. FLGROSS      0.18060    0.01895   9.532 < 2e-16 ***
17. FPOLL       -0.03229    0.01352  -2.388  0.01770 *
18. FLTOTMED    -0.02719    0.01298  -2.095  0.03725 *
19. F03H24      0.05688    0.02537   2.242  0.02585 *
20. FTEH24     -0.06978    0.02516  -2.774  0.00598 **
21. FSATEM      0.03480    0.01580   2.202  0.02860 *
22. FLGEW       0.10461    0.01915   5.463 1.18e-07 ***
23. FSPFEI      0.02331    0.01541   1.513  0.13172
24. ---
25. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
26.
27. Residual standard error: 0.2001 on 239 degrees of freedom
28. Multiple R-squared:  0.6957,    Adjusted R-squared:  0.683
29. F-statistic: 54.65 on 10 and 239 DF,  p-value: < 2.2e-16
30.
31. > summary(forward_AIC)
32.
33. Call:
34. lm(formula = train_set$FFVC ~ FLGROSS + SEX + FLGEW + FSATEM +
35.     FLTOTMED + FPOLL + AMATOP + FTEH24 + F03H24 + FSPFEI, data = train_set)
36.
37. Residuals:
38.      Min      1Q   Median      3Q      Max
39. -0.47336 -0.13900 -0.00203  0.13173  0.55450
40.
41. Coefficients:
42.      Estimate Std. Error t value Pr(>|t|)
43. (Intercept)  2.28480    0.01266 180.534 < 2e-16 ***
44. FLGROSS      0.18060    0.01895   9.532 < 2e-16 ***
45. SEX         -0.09756    0.01321  -7.386 2.51e-12 ***
46. FLGEW       0.10461    0.01915   5.463 1.18e-07 ***
47. FSATEM      0.03480    0.01580   2.202  0.02860 *
48. FLTOTMED    -0.02719    0.01298  -2.095  0.03725 *
49. FPOLL       -0.03229    0.01352  -2.388  0.01770 *
50. AMATOP      -0.02142    0.01324  -1.618  0.10703
51. FTEH24     -0.06978    0.02516  -2.774  0.00598 **

```



```

52. F03H24      0.05688    0.02537    2.242  0.02585 *
53. FSPFEI      0.02331    0.01541    1.513  0.13172
54. ---
55. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
56.
57. Residual standard error: 0.2001 on 239 degrees of freedom
58. Multiple R-squared:  0.6957,    Adjusted R-squared:  0.683
59. F-statistic: 54.65 on 10 and 239 DF,  p-value: < 2.2e-16
60.
61. > summary(backward_BIC)
62.
63. Call:
64. lm(formula = FFVC ~ SEX + FLGROSS + FPOLL + FLTOTMED + F03H24 +
65.     FTEH24 + FSATEM + FLGEW, data = train_set)
66.
67. Residuals:
68.      Min       1Q   Median       3Q      Max
69. -0.46964 -0.14614  0.00336  0.13200  0.56232
70.
71. Coefficients:
72.             Estimate Std. Error t value Pr(>|t|)
73. (Intercept)  2.28480    0.01273 179.537 < 2e-16 ***
74. SEX          -0.09755    0.01320  -7.390 2.40e-12 ***
75. FLGROSS       0.17628    0.01892   9.318 < 2e-16 ***
76. FPOLL        -0.03219    0.01355  -2.375 0.018312 *
77. FLTOTMED     -0.02928    0.01297  -2.257 0.024919 *
78. F03H24       0.06308    0.02529   2.494 0.013291 *
79. FTEH24      -0.07099    0.02520  -2.817 0.005243 **
80. FSATEM       0.04508    0.01338   3.369 0.000878 ***
81. FLGEW       0.10742    0.01900   5.653 4.44e-08 ***
82. ---
83. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
84.
85. Residual standard error: 0.2012 on 241 degrees of freedom
86. Multiple R-squared:  0.6898,    Adjusted R-squared:  0.6795
87. F-statistic: 66.98 on 8 and 241 DF,  p-value: < 2.2e-16
88.
89. > summary(forward_BIC)
90.
91. Call:
92. lm(formula = train_set$FFVC ~ FLGROSS + SEX + FLGEW + FSATEM +
93.     FLTOTMED + FPOLL, data = train_set)
94.
95. Residuals:

```

```

96.      Min      1Q   Median      3Q      Max
97. -0.49700 -0.13431 -0.00183  0.12576  0.53832
98.
99. Coefficients:
100.      Estimate Std. Error t value Pr(>|t|)
101. (Intercept)  2.28480     0.01288 177.375 < 2e-16 ***
102. FLGROSS      0.17377     0.01912   9.091 < 2e-16 ***
103. SEX         -0.09719     0.01334  -7.287 4.42e-12 ***
104. FLGEW       0.10694     0.01922   5.565 6.91e-08 ***
105. FSATEM      0.04752     0.01346   3.531 0.000495 ***
106. FLTOTMED    -0.03165     0.01310  -2.416 0.016423 *
107. FPOLL       -0.03201     0.01366  -2.343 0.019957 *
108. ---
109. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
110.
111. Residual standard error: 0.2037 on 243 degrees of freedom
112. Multiple R-squared:  0.6795,    Adjusted R-squared:  0.6716
113. F-statistic: 85.87 on 6 and 243 DF,  p-value: < 2.2e-16

```

PREDICTIONS ON THEIR PERFORMANCE:

In my opinion, the reduced model shall perform better than the full model.

At first, it is easy to find out that AIC/BIC of these models shall be lower than the full model, which indicates that the reduced models are better. And besides, according to the Std. Error of our estimates, nearly all estimates of the reduced model are more stable than the full model. In addition, notice that the p -values of our reduced models show that nearly all estimates are significantly not 0, while the p -values of the full model do not. This may result in an overfitting in the full model.

Thus, we could conclude that the reduced models are supposed to be better than the full model, i.e. has a lower prediction error.

Task 4

CODE:

```

1. #4
2. #select lasso
3. ##split data
4. library(glmnet)
5. n_random = sample(1:496,496)
6. train_set <- data[c(sort(n_random[1:250])),]
7. X <- train_set[-25]
8. y <- train_set[25]
9.

```

```

10. ##setting range of parameters
11. lambdas <- seq(0.0001, 0.2, by = 0.01)
12. n_lambdas = length(lambdas)
13.
14. ##5 fold cv
15. err_cv = numeric(n_lambdas)
16. for (k in 1:n_lambdas){
17.   n_random_cv = sample(1:250,250)
18.   n_cv = c(sort(n_random_cv[1:50]), sort(n_random_cv[51:100]),
19.             sort(n_random_cv[101:150]), sort(n_random_cv[151:200]),
20.             sort(n_random_cv[201:250]))
21.   for (i in 1:5){
22.     a <- 50*i-49
23.     b <- 50*i
24.     X_train <- data.matrix(X[-n_cv[a:b], ])
25.     y_train <- data.matrix(y[-n_cv[a:b], ])
26.     X_test <- data.matrix(X[n_cv[a:b], ])
27.     y_test <- data.matrix(y[n_cv[a:b], ])
28.     lasso_reg <- glmnet(X_train, y_train, alpha = 1, lambda = lambdas[k])
29.     err_cv[k] <- err_cv[k] + sum((predict(lasso_reg, newx = X_test , s = lambdas[k]
    ) - y_test)^2)/50
30.   }
31.   err_cv[k] <- err_cv[k]/5
32. }
33.
34. ##bootstrap
35. err_boots = numeric(n_lambdas)
36. for (k in 1:n_lambdas){
37.   for (i in 1:10){
38.     n_random_boots = sample(1:250,250, replace = TRUE)
39.     n_boots = c(sort(n_random_boots))
40.     X_train <- data.matrix(X[-n_boots, ])
41.     y_train <- data.matrix(y[-n_boots, ])
42.     X_test <- data.matrix(X[n_boots, ])
43.     y_test <- data.matrix(y[n_boots, ])
44.     lasso_reg <- glmnet(X_train, y_train, alpha = 1, lambda = lambdas[k])
45.     err_boots[k] <- err_boots[k] + sum((predict(lasso_reg, newx = X_test , s = lamb
    das[k]) - y_test)^2)/length(y_test)
46.   }
47.   err_boots[k] <- err_boots[k]/10
48. }
49.
50. ##plot
51. library(ggplot2)

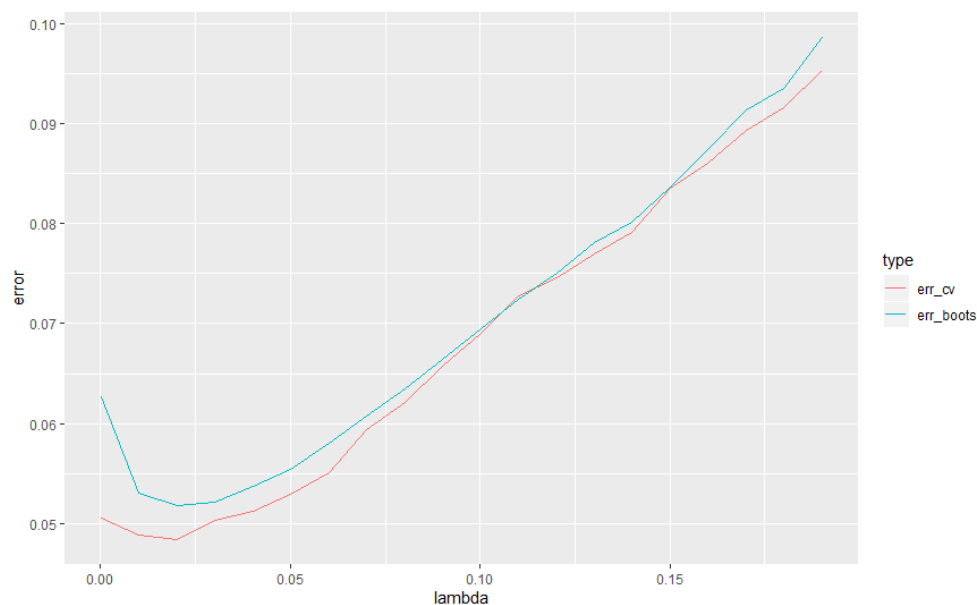
```

```

52. df1 <- data.frame(lambda = lambdas, error = err_cv, type = "err_cv")
53. df2 <- data.frame(lambda = lambdas, error = err_boots, type = "err_boots")
54. df = rbind(df1, df2)
55. ggplot(df) + geom_line(aes(lambda, error, colour=type))

```

RESULT:



COMMENTS:

We see that, for Bootstrap procedure, the optimal lambda is around 0.02, while the optimal lambda for 5-fold cross-validation is also around 0.02.

By repeating the test several times (not shown here), although the result of Bootstrap procedure could vary, it shall still locate in (0.01, 0.03), while the result of 5-fold CV also has a similar bound, i.e. (0.01, 0.03).

Thus, we could conclude that the results of these 2 procedures agree with each other.

Task 5

Obviously, those categorical variables are not supposed to contain non-linear effects, since they are actually binary variables, i.e. 1 or 0.

Although we only need to consider the numerical variables, we are not knowing which variables contain non-linear effects. And thus, we just fit a GAM with splines for all numerical variables. Afterwards, we could check the ANOVA and the plots to justify which variables may contain non-linear effects.

We are using **mgcv** package to implement the GAM model, since such a package is able to automatically select the smoothness parameter, it shall use $GCV-C_p$ by default, but other methods, such

as AIC and etc., are also available.

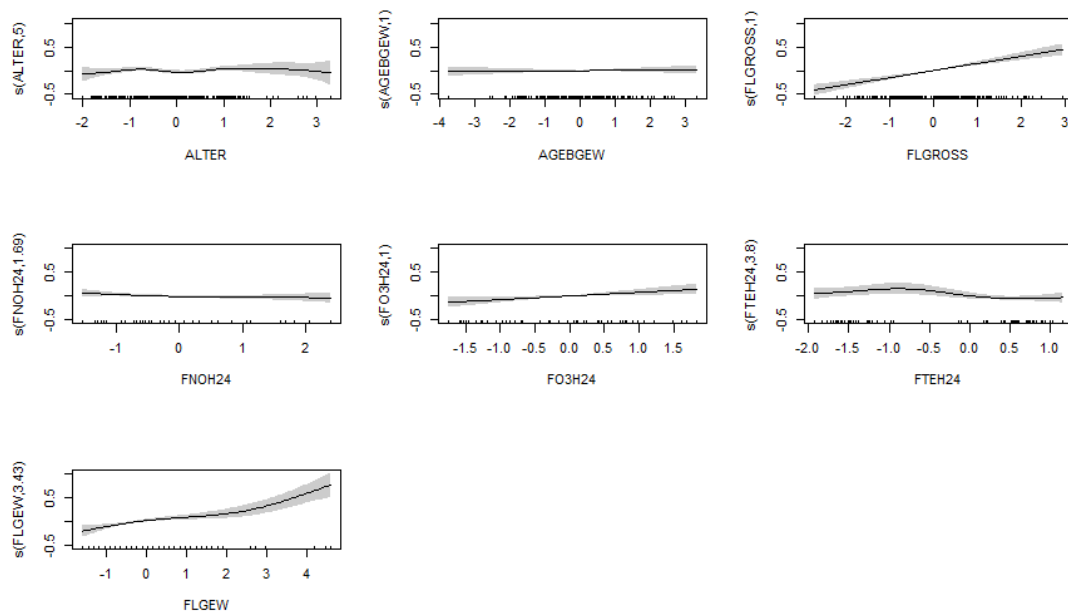
CODE:

```
1. #5
2. ##GAM
3. library(mgcv)
4. model_gam_test <- gam(train_set$FFVC~s(ALTER)+s(AGEBGEW)+s(FLGROSS)+s(FNOH24)+s(F03
  H24)+s(FTEH24)+s(FLGEW)+ADHEU+SEX+HOCHOZON+AMATOP+AVATOP+ADEKZ+ARAUCH+FSNIGHT+FMILB
  +FTIER+FPOLL+FLTOTMED+FSPT+FSATEM+FSAUGE+FSPFEI+FSHLAUF, data = train_set)
5. summary(model_gam_test)
6. plot(model_gam_test, pages = 1, scheme = 1, all.terms = TRUE)
```

RESULT:

```
1. Approximate significance of smooth terms:
2.           edf Ref.df      F  p-value
3. s(ALTER)   5.004  6.099  1.191   0.3323
4. s(AGEBGEW) 1.000  1.000  0.146   0.7027
5. s(FLGROSS) 1.000  1.000 52.219 6.34e-12 ***
6. s(FNOH24)  1.687  2.100  1.440   0.2493
7. s(F03H24)  1.000  1.000  6.025   0.0149 *
8. s(FTEH24)  3.797  4.664  1.812   0.1192
9. s(FLGEW)   3.433  4.259 10.098 1.03e-07 ***
10. ---
11. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
12.
13. R-sq.(adj) =  0.715   Deviance explained = 75.3%
14. GCV = 0.04159   Scale est. = 0.03578    n = 250
```

PLOT:



By analyzing the “Approximate significance of smooth terms” and the plots above, it is clear that only the **FLGEW** may have a significant non-linear effect in such case. We then just implement a GAM which only applies splines to **FLGEW**.

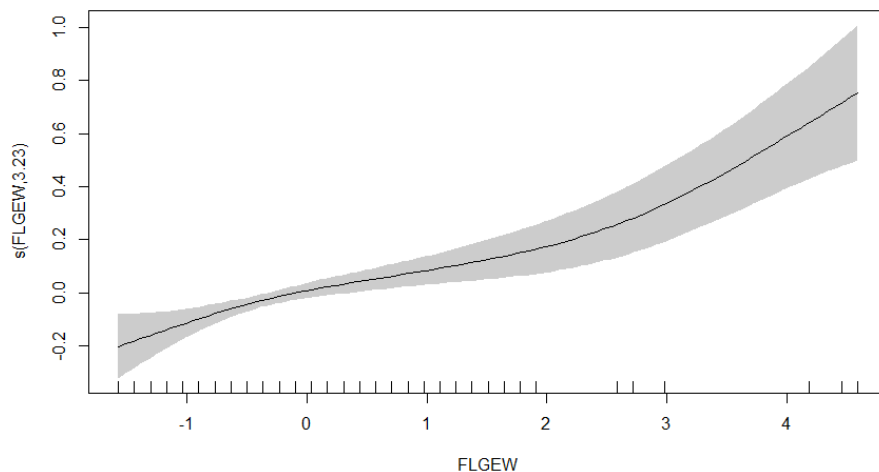
CODE:

```
1. model_gam <- gam(train_set$FFVC~ALTER+AGEBGEW+FLGROSS+FNOH24+FO3H24+FTEH24+s(FLGEW)
+ADHEU+SEX+HOCHOZON+AMATOP+AVATOP+ADEKZ+ARAUCH+FSNIGHT+FMILB+FTIER+FPOLL+FLTOTMED+
FSPT+FSATEM+FSAUGE+FSPFEI+FSHLAUF, data = train_set)
2. summary(model_gam)
3. plot(model_gam, pages = 1, scheme = 1, all.terms = FALSE)
```

RESULT:

```
1. Approximate significance of smooth terms:
2.          edf Ref.df      F  p-value
3. s(FLGEW) 3.226  4.028 10.54 6.85e-08 ***
4. ---
5. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
6.
7. R-sq.(adj) =  0.699   Deviance explained = 73.1%
8. GCV = 0.042305   Scale est. = 0.037698   n = 250
```

PLOT:



COMMENTS:

The statistics and the plot above are the same as we expected. The p -value shows that our smooth terms are significantly not 0, and besides, the adjusted R-square is much higher than the original linear model which we implemented before. It means that this model fits better to the data, and therefore the non-linearity could be captured.

As for the plot above, it is clearly shown that the non-linearity of **FLGEW** does exist. And we then just try to implement a linear model with adding polynomial terms to this variable. We add the terms from degree 3 and keep increasing the degree until the ANOVA result is not proper any more.

CODE:

```
1. ##GAM polynomial
2. model_pol <- lm(train_set$FFVC ~ I(FLGEW^2) + I(FLGEW^3) + ., data = train_set)
3. summary(model_pol)
4. model_pol <- lm(train_set$FFVC ~ I(FLGEW^2) + I(FLGEW^3) + I(FLGEW^4) + ., data = train_set)
5. summary(model_pol)
```

RESULT:

```
1. Call:
2. lm(formula = train_set$FFVC ~ I(FLGEW^2) + I(FLGEW^3) + ., data = train_set)
3.
4. Residuals:
5.      Min       1Q   Median       3Q      Max
6. -0.5054 -0.1296  0.0074  0.1166  0.4875
7.
8. Coefficients:
9.              Estimate Std. Error t value Pr(>|t|)
```

```

10. (Intercept)  2.3029252  0.0173051 133.078 < 2e-16 ***
11. I(FLGEW^2)  -0.0539251  0.0164036  -3.287  0.00117 **
12. I(FLGEW^3)   0.0118474  0.0040024   2.960  0.00341 **
13. ALTER        0.0192996  0.0146568   1.317  0.18927
14. ADHEU        0.0002489  0.0139112   0.018  0.98574
15. SEX         -0.1018458  0.0127167  -8.009  6.38e-14 ***
16. HOCHOZON    -0.0367559  0.0180419  -2.037  0.04281 *
17. AMATOP       0.0172865  0.0138462   1.248  0.21317
18. AVATOP      -0.0118461  0.0131629  -0.900  0.36911
19. ADEKZ       -0.0162862  0.0134587  -1.210  0.22752
20. ARAUCH      -0.0076738  0.0130251  -0.589  0.55635
21. AGEBGEW     0.0051167  0.0136240   0.376  0.70760
22. FSNIGHT     -0.0129540  0.0148930  -0.870  0.38534
23. FLGROSS     0.1473639  0.0214269   6.878  6.06e-11 ***
24. FMILB      -0.0083075  0.0173360  -0.479  0.63226
25. FNOH24     -0.0150290  0.0172761  -0.870  0.38527
26. FTIER       0.0119805  0.0144093   0.831  0.40661
27. FPOLL      -0.0232736  0.0248676  -0.936  0.35034
28. FLTOTMED   -0.0122225  0.0133976  -0.912  0.36260
29. FO3H24     0.0205912  0.0287983   0.715  0.47535
30. FSPT       -0.0046885  0.0292319  -0.160  0.87272
31. FTEH24     -0.0045167  0.0264904  -0.171  0.86477
32. FSATEM     0.0425879  0.0169139   2.518  0.01251 *
33. FSAUGE     0.0046250  0.0139578   0.331  0.74069
34. FLGEW      0.1054947  0.0239645   4.402  1.66e-05 ***
35. FSPFEI     0.0205858  0.0176614   1.166  0.24503
36. FSHLAUF    -0.0261818  0.0148196  -1.767  0.07865 .
37. ---
38. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
39.
40. Residual standard error: 0.1919 on 223 degrees of freedom
41. Multiple R-squared:  0.7213,    Adjusted R-squared:  0.6888
42. F-statistic: 22.2 on 26 and 223 DF,  p-value: < 2.2e-16
43.
44. Call:
45. lm(formula = train_set$FFVC ~ I(FLGEW^2) + I(FLGEW^3) + I(FLGEW^4) +
46.     ., data = train_set)
47.
48. Residuals:
49.      Min       1Q   Median       3Q      Max
50. -0.51798 -0.12306  0.00294  0.11571  0.47478
51.
52. Coefficients:
53.             Estimate Std. Error t value Pr(>|t|)

```



```

54. (Intercept)  2.3011418  0.0172555 133.357 < 2e-16 ***
55. I(FLGEW^2)  -0.0461002  0.0169284  -2.723  0.00698 **
56. I(FLGEW^3)  -0.0119088  0.0141377  -0.842  0.40050
57. I(FLGEW^4)   0.0043722  0.0024965   1.751  0.08127 .
58. ALTER       0.0182333  0.0146020   1.249  0.21310
59. ADHEU       -0.0005823  0.0138553  -0.042  0.96652
60. SEX         -0.1032830  0.0126847  -8.142  2.78e-14 ***
61. HOCHOZON    -0.0385558  0.0179883  -2.143  0.03317 *
62. AMATOP      0.0182188  0.0137927   1.321  0.18789
63. AVATOP     -0.0125921  0.0131092  -0.961  0.33782
64. ADEKZ      -0.0157439  0.0134003  -1.175  0.24129
65. ARAUCH     -0.0055920  0.0130195  -0.430  0.66797
66. AGEBGEW     0.0043560  0.0135683   0.321  0.74848
67. FSNIGHT    -0.0145044  0.0148509  -0.977  0.32980
68. FLGROSS     0.1473226  0.0213283   6.907  5.14e-11 ***
69. FMILB      -0.0032768  0.0174936  -0.187  0.85159
70. FNOH24     -0.0163209  0.0172124  -0.948  0.34406
71. FTIER       0.0089142  0.0144494   0.617  0.53792
72. FPOLL      -0.0233969  0.0247532  -0.945  0.34558
73. FLTOTMED   -0.0108870  0.0133578  -0.815  0.41593
74. FO3H24     0.0219318  0.0286759   0.765  0.44519
75. FSPT       -0.0042973  0.0290982  -0.148  0.88273
76. FTEH24     -0.0041095  0.0263695  -0.156  0.87630
77. FSATEM     0.0408472  0.0168654   2.422  0.01624 *
78. FSAUGE     0.0047144  0.0138936   0.339  0.73469
79. FLGEW       0.1457747  0.0331363   4.399  1.69e-05 ***
80. FSPFEI     0.0220139  0.0175990   1.251  0.21230
81. FSHLAUF    -0.0286810  0.0148203  -1.935  0.05423 .
82. ---
83. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
84.
85. Residual standard error: 0.191 on 222 degrees of freedom
86. Multiple R-squared:  0.7251,    Adjusted R-squared:  0.6917
87. F-statistic: 21.69 on 27 and 222 DF,  p-value: < 2.2e-16

```

COMMENTS:

According to the results above, we see that, for a polynomial of degree 3, both **FLGEW**² and **FLGEW**³ are significantly not 0. And besides, recall the R^2 for the original linear model we implemented before, we find that our new model has a much higher R^2 , i.e. a better fitness.

Now, we could conclude that our new models really capture some of the non-linearity, and both GAM with splines and the GLM with polynomial terms are supposed to perform well. But as for the further comparisons, we could not guarantee anything more until we really test the performance of our models, i.e. Task 7.

Task 6

FYI, since all our data are standardized before fitting, we are then using “center = TRUE” for glmboost, i.e. **linear model** for component-wise boosting. And according to what we have stated before, only the numerical variables shall be applied a component-wise boosting with **splines** as base learners. And since the response of our data, i.e. FFVC, is numerical, so we are using default Gaussian family for all following functions:

CODE:

```
1. library(mboost)
2. compo_boosting_lin <- glmboost(FFVC ~ ., data = train_set, center = TRUE)
3. compo_boosting_spl <- mboost(FFVC ~ ALTER+AGEBGEW+FLGROSS+FNOH24+F03H24+FTEH24+FLGEW, data = train_set, baselearner = 'bbs')
4. compo_boosting_tree <- mboost(FFVC ~ ., data = train_set, baselearner = "btree")
5. summary(compo_boosting_lin)
6. summary(compo_boosting_spl)
7. summary(compo_boosting_tree)
```

Thus, we are able to produce the **variable selection frequencies**:

Linear Models:

```
1. Selection frequencies:
2. SEX  FLGROSS  FLGEW  FLTOTMED  FSATEM  FSPFEI  FPOLL  ADHEU  ALTER  FNOH24  HOCHOZON  FMILB
3. 0.21   0.17   0.12   0.10   0.08   0.08   0.07   0.05   0.04   0.04   0.03   0.01
```

Splines:

```
1. Selection frequencies:
2. bbs(FTEH24) bbs(FLGROSS)  bbs(FLGEW)  bbs(ALTER)  bbs(FNOH24) bbs(AGEBGEW)  bbs(F03H24)
3.      0.26      0.21      0.15      0.14      0.12      0.09      0.03
```

Trees:

```
1. Selection frequencies:
2. btree(FLGROSS)btree(FLGEW) btree(SEX) btree(ALTER) btree(AGEBGEW)btree(HOCHOZON) btree(F03H24)
3.      0.30      0.20      0.15      0.08      0.06      0.05      0.05
4. btree(FNOH24)  btree(ADHEU)  btree(FTEH24)  btree(FSATEM)  btree(FPOLL)  btree(FLTOTMED)
5.      0.03      0.02      0.02      0.02      0.01      0.01
```

And the **regression coefficients** of linear models:

```
1. Coefficients:
2. (Intercept)      ALTER      ADHEU      SEX      HOCHOZON
3. -4.2850504427  0.0100680434 -0.0333030854 -0.1967860125  0.0080105588
```

4.	FLGROSS	FMILB	FNOH24		
5.	0.0291114253	-0.0032612557	-0.0002388067		
6.	FPOLL	FLTOTMED	FSATEM	FLGEW	FSPFEI
7.	-0.0268506402	-0.0418009858	0.0822904467	0.0097438394	0.0613377048

Task 7

CODE:

```

1. #7
2. ##train error
3. train_error <- NULL
4. train_error['full'] <- mean(residuals(model)^2)
5.
6. train_error['backAIC'] <- mean(residuals(backward_AIC)^2)
7. train_error['forAIC'] <- mean(residuals(forward_AIC)^2)
8. train_error['backBIC'] <- mean(residuals(backward_BIC)^2)
9. train_error['forBIC'] <- mean(residuals(forward_BIC)^2)
10.
11. train_error['lasso_cv'] <- mean((data.matrix(y_task4) - predict(lasso_cv, newx = data.matrix(X_task4), s = lambda_cv))^2)
12. train_error['lasso_boots'] <- mean((data.matrix(y_task4) - predict(lasso_boots, newx = data.matrix(X_task4), s = lambda_boots))^2)
13.
14. train_error['gam'] <- mean(residuals(model_gam)^2)
15. train_error['pol'] <- mean(residuals(model_pol)^2)
16.
17. train_error['com_boost_lin'] <- mean(residuals(compo_boosting_lin)^2)
18. train_error['com_boost_spl'] <- mean(residuals(compo_boosting_spl)^2)
19. train_error['com_boost_tree'] <- mean(residuals(compo_boosting_tree)^2)
20.
21. ##test error
22. X_task7 <- test_set[-25]
23. y_task7 <- data.matrix(test_set[25])
24. test_error <- NULL
25.
26. test_error['full'] <- mean((y_task7 - predict(model, newdata = X_task7))^2)
27. test_error['backAIC'] <- mean((y_task7 - predict(backward_AIC, newdata = X_task7))^2)
28. test_error['forAIC'] <- mean((y_task7 - predict(forward_AIC, newdata = X_task7))^2)
29. test_error['backBIC'] <- mean((y_task7 - predict(backward_BIC, newdata = X_task7))^2)
30. test_error['forBIC'] <- mean((y_task7 - predict(forward_BIC, newdata = X_task7))^2)
31.

```

```

32. test_error['lasso_cv'] <- mean((y_task7 - predict(lasso_cv, newx = data.matrix(X_task7), s = lambda_cv))^2)
33. test_error['lasso_boots'] <- mean((y_task7 - predict(lasso_boots, newx = data.matrix(X_task7), s = lambda_boots))^2)
34.
35. test_error['gam'] <- mean((y_task7 - predict(model_gam, newdata = X_task7))^2)
36. test_error['pol'] <- mean((y_task7 - predict(model_pol, newdata = X_task7))^2)
37.
38. test_error['com_boost_lin'] <- mean((y_task7 - predict(compo_boosting_lin, newdata = X_task7))^2)
39. test_error['com_boost_spl'] <- mean((y_task7 - predict(compo_boosting_spl, newdata = X_task7[num_vari]))^2)
40. test_error['com_boost_tree'] <- mean((y_task7 - predict(compo_boosting_tree, newdata = X_task7))^2)
41.
42. print("train error")
43. print(train_error)
44. print("test error")
45. print(test_error)

```

RESULT:

```

1. [1] "train error"
2.   full      backAIC      forAIC      backBIC      forBIC
3. 0.03950484 0.04123226 0.04123226 0.04123226 0.04123226
4.  lasso_cv  lasso_boots      gam      pol      com_boost_lin
5. 0.10981893 0.11133505 0.03864648 0.03879739 0.04236932
6. com_boost_spl com_boost_tree
7. 0.04686207 0.03688197
8.
9. [1] "test error"
10.  full      backAIC      forAIC      backBIC      forBIC
11. 0.04831460 0.04789635 0.04789635 0.04789635 0.04789635
12.  lasso_cv  lasso_boots      gam      pol      com_boost_lin
13. 0.14082055 0.14424095 0.04735161 0.04607857 0.04687056
14. com_boost_spl com_boost_tree
15. 0.05514283 0.05273946

```

COMMENTS:

We see overall that the training error are smaller than test error, which is exactly the same as we expected, since we are accessing the training set when we are fitting our model, while the test set keeps unseen.

We notice that both the training error and the prediction error for lasso regression are much larger

than others. This could be a result of standardizing. We have standardized both numerical variables and categorical variables into mean= 0 and sd.= 1, and it means that the corresponding **estimates** of these variables could **be shrunk/zoomed**, and this shall influence the LASSO.

Recall our prediction in **Task 3**, we see that the reduced models have larger training errors, but lower test errors. It is the same as we predicted. Although the full model fits the training data better, at the same time, it also brings more risk of overfitting. But as for the reduced models, they did not fit the training data as well as the full model, but they could take uncertainty of new data into account and avoid the overfitting issues, and therefore, they have better performance when being tested.

Problem 2

Before all tasks: split the data

By using the same method in **Problem 1**, we at first check the distribution for diabetes (positive/negative)

CODE:

```
1. library(mlbench)
2. library(kknn)
3. data(PimaIndiansDiabetes)
4. dat <- PimaIndiansDiabetes
5.
6. a0 <- nrow(dat[dat['diabetes'] == "pos",])
7. a1 <- nrow(dat[dat['diabetes'] == "neg",])
8. print(c("pos", "neg"))
9. print(c(a0, a1))
```

RESULT:

```
1. [1] "pos" "neg"
2. [1] 268 500
```

It seems that we could just divide the data randomly, since there is no extreme distribution here.

Thus, we just split the train set and test set randomly and check the distribution of **diabetes** again to ensure that the training set and test set have similar distributions.

CODE:

```
1. a0 <- nrow(train_set[train_set['diabetes'] == "pos",])
2. a1 <- nrow(train_set[train_set['diabetes'] == "neg",])
```

```

3. print(c("pos", "neg"))
4. print(c(a0, a1))
5.
6. a0 <- nrow(test_set[test_set['diabetes'] == "pos",])
7. a1 <- nrow(test_set[test_set['diabetes'] == "neg",])
8. print(c("pos", "neg"))
9. print(c(a0, a1))

```

RESULT:

```

1. [1] "pos" "neg"
2. [1] 183 327
3.
4. [1] "pos" "neg"
5. [1] 85 173

```

Obviously, our train set and test set follow similar distributions of **diabetes**.

Task 1

CODE:

```

1. #task 1
2. ##5-fold cv
3. err_cv <- numeric(k_range)
4. n_random_cv = sample(1:510, 510)
5. n = c(sort(n_random_cv[1:102]), sort(n_random_cv[103:204]), sort(n_random_cv[205:306]),
6.       sort(n_random_cv[307:408]), sort(n_random_cv[409:510]))
7. k_range = 50
8.
9. for (i in 1:k_range){
10.   for(j in 1:5){
11.     a <- 102*j-101
12.     b <- 102*j
13.     fit <- kknn(diabetes ~ ., train = train_set[-
14.               n[a:b], ], test = train_set[n[a:b], ], k = i)
15.     err_cv[i] <- err_cv[i] + sum(fitted.values(fit)!=train_set[n[a:b],9])/102
16.   }
17.   err_cv[i] <- err_cv[i]/5
18. }
19. ##LOOCV

```

```

20. err_LOOCV <- numeric(k_range)
21. for (i in 1:k_range){
22.   for(j in 1:510){
23.     fit <- kknn(diabetes ~ ., train = train_set[-
      j, ], test = train_set[j, ], k = i)
24.     err_LOOCV[i] <- err_LOOCV[i] + sum(fitted.values(fit)!=train_set[j, 9])
25.   }
26.   err_LOOCV[i] <- err_LOOCV[i]/510
27. }
28.
29. ##test
30. err_test <- numeric(k_range)
31. for (i in 1:k_range){
32.   fit <- kknn(diabetes ~ ., train = train_set, test = test_set, k = i)
33.   err_test[i] <- sum(fitted.values(fit)!=test_set[,9])/258
34. }
35.
36. ##print&plot
37. which.min(err_cv)
38. which.min(err_LOOCV)
39. which.min(err_test)
40.
41. plot(1:k_range, err_cv, ylim=c(0.2, 0.4), xlab="k", ylab="error",col="blue", type="
  o",lwd=2)
42. lines(1:k_range, err_LOOCV, col="green",type="o",lwd=2)
43. lines(1:k_range, err_test, col="black",type="o",lwd=2)
44. legend("topleft", c("error_cv","error_LOOCV","error_test"),fill=c("blue","green","b
  lack"))

```

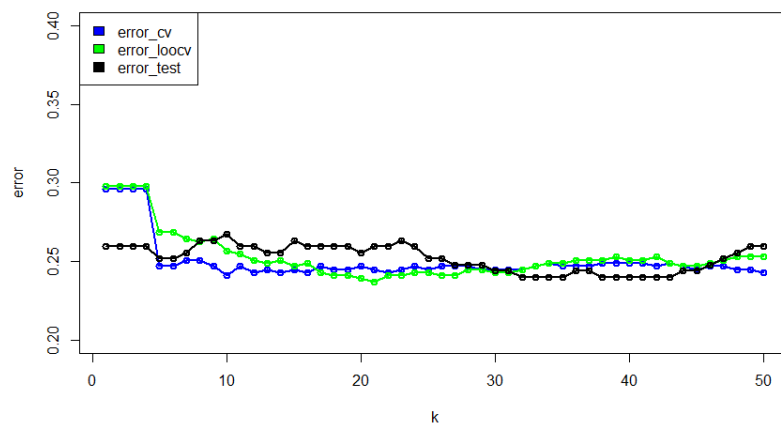
RESULT:

```

1. > which.min(err_cv)
2. [1] 24
3. > which.min(err_LOOCV)
4. [1] 24
5. > which.min(err_test)
6. [1] 19

```

PLOT:



COMMENTS:

Obviously, we see that all the 3 types of errors have similar tendencies, and it shows that both 5-fold CV and LOOCV could represent the distribution of the whole data set to some degree.

Task 2

By using the package **mgcv** with setting “select = TRUE”, we could penalize some of the variables into 0. Such a procedure is pretty much similar with a LASSO regression. And it could automatically process the variable selection.

CODE:

```
1. model_gam <- mgcv::gam(as.numeric(diabetes)~s(pregnant)+s(glucose)+s(pressure)+s(triceps)+s(insulin)+s(mass)+s(pedigree)+s(age), data = train_set, select = TRUE)
2. summary(model_gam)
3. plot(model_gam, pages = 1, scheme = 1, all.terms = FALSE)
```

RESULT:

```
1. Approximate significance of smooth terms:
2.      edf Ref.df      F p-value
3. s(pregnant) 4.863e+00      9  1.534 0.009605 **
4. s(glucose)  2.463e+00      9 14.331 < 2e-16 ***
5. s(pressure) 9.717e-01      9  1.270 0.000587 ***
6. s(triceps)  4.408e-08      9  0.000 0.851817
7. s(insulin)  8.560e-08      9  0.000 0.604047
8. s(mass)     1.241e+00      9  2.272 2.07e-06 ***
9. s(pedigree) 9.033e-01      9  0.337 0.039042 *
10. s(age)     6.087e+00      9  3.056 5.93e-05 ***
11. ---
```

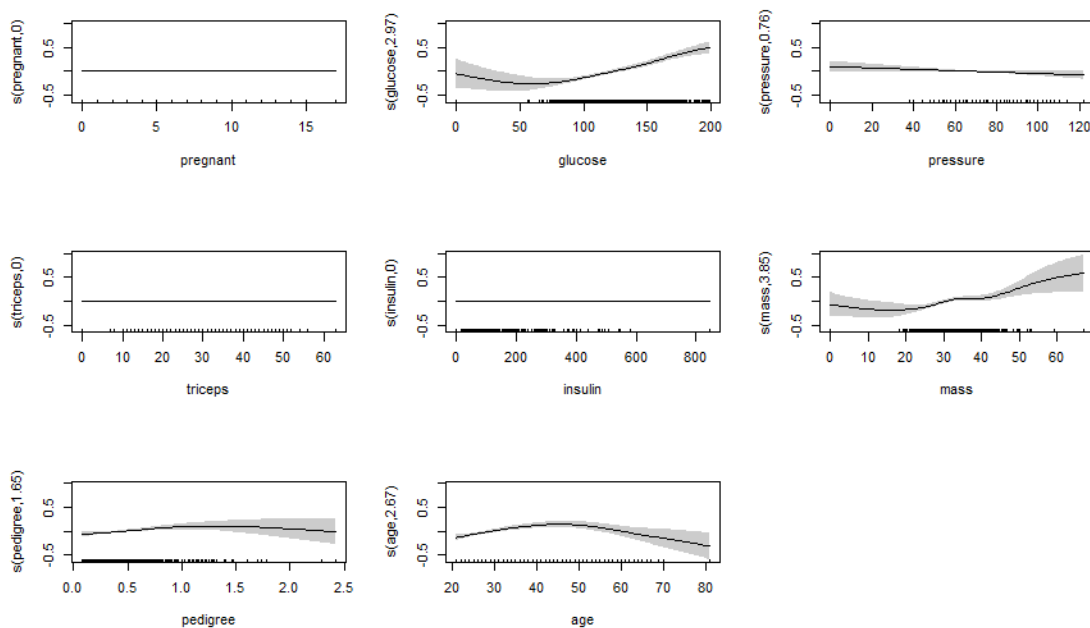


```

12. Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
13.
14. R-sq.(adj) =  0.395   Deviance explained = 41.5%
15. GCV = 0.14292   Scale est. = 0.13801   n = 510

```

PLOT:



According to the result and plot above, we see that some variables are reduced into 0, which is equivalent with removing these variables. Besides, all the **Effective Degrees of Freedom** are also automatically decided by penalizing.

And thus, this model could now be considered as the best model after processing subset selection.

Task 3

TREE:

In order to generate a proper classification tree, we need to generate a huge tree at first and then try to prune it into the best one.

CODE:

```

1. library(tree)
2. n_random = sample(1:510,510)
3. tree_train_set <- dat[c(sort(n_random[1:340])),]
4. tree_test_set <- dat[c(sort(n_random[341:510])),]
5.

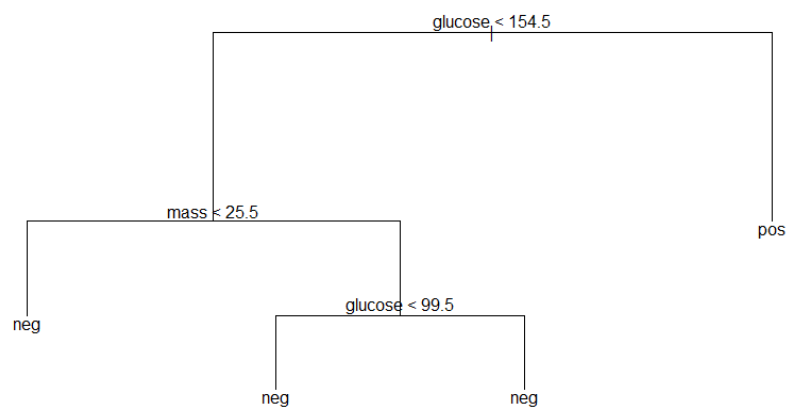
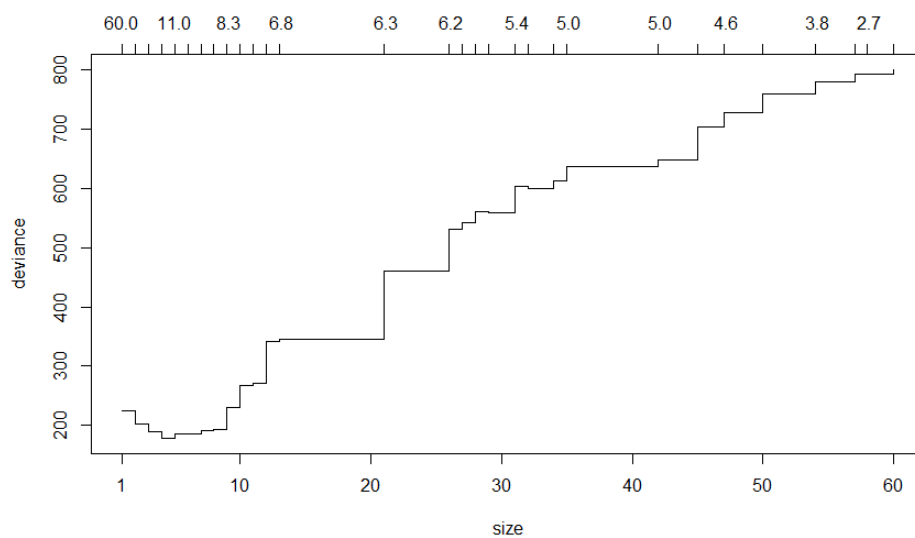
```

```

6. t1<- tree(diabetes~., data = tree_train_set, control = tree.control(nobs = 340, min
   size = 0, mindev = 0))
7. plot(t1)
8. text(t1)
9.
10. trees<- prune.tree(t1, newdata = tree_test_set)
11. plot(trees)
12.
13. tree_best <- prune.tree(t1, best = trees$size[trees$dev == min(trees$dev)])
14. plot(tree_best)
15. text(tree_best)

```

RESULT:



BAGGING:

By directly using the function **bagging** in the package **ipred**, we could simply implement a bagging trees model. In such case, there is not clearly improvement even if we increase the number of baggings, i.e. **nbagg**. Thus, we just use an empirical number 100.

CODE:

```
1. library(ipred)
2. bagging(diabetes~., nbagg=100, method="standard", data = train_set)
```

RESULT:

```
1. Bagging classification trees with 100 bootstrap replications
2.
3. Call: bagging.data.frame(formula = diabetes ~ ., data = train_set,
4.     nbagg = 100, method = "standard", coob = TRUE)
5.
6. Out-of-bag estimate of misclassification error: 0.2588
```

RANDOM FOREST:

According to the lectures, there might be pretty little improvement with tuning parameters, thus, we just directly implement the model by using package **randomForest** with default parameters. But notice that we need to at least ensure the parameter **ntree** is not too small, i.e. 1, 2 and etc., in order to generate enough amount of trees.

We just simply use an empirical number **ntree= 25**.

CODE:

```
1. library(randomForest)
2. randomForest(diabetes~., ntree = 25, data = train_set)
```

RESULT:

```
1.          OOB estimate of  error rate: 26.08%
2. Confusion matrix:
3.      neg pos  class.error
4. neg 271  57   0.1737805
5. pos  76 106   0.4175824
```

NEURAL NETWORK:

By using package **neuralnet**, we could simply implement a neural network model for our data. Since the package is using **rprop+**, so we do not need to tune the learning rate any more, and instead, it could be automatically adjusted.

CODE:

```
1. library(neuralnet)
2. model_nn <- neuralnet(diabetes~., data = train_set)
```

RESULT:

The details of neural networks are too much to be shown here.

ADABOOST:

In fact, this algorithm is pretty simple to be implemented, but in order to make the model more universal, i.e. could be directly used with **predict** and etc., we just use the package **fastAdaboost**, and since this package uses C++, as it is called, it runs faster, so we could just choose a larger amount of classifiers, **nIter= 200**.

CODE:

```
1. library(fastAdaboost)
2. model_ada <- adaboost(diabetes~., data = train_set, nIter = 200)
```

Task 4

Anyway, the numbers are reliable, thus, we just use the test set to test the performance at first.

CODE:

```
1. #task4
2. ###test
3. error_tree <- sum((predict(model_tree, newdata = test_set[, -9])[,2]>0.5)!=
4.                  (test_set[,9]=="pos"))/258
5. error_bagging <- sum(predict(model_bagging, newdata = test_set[, -9])!=
6.                      test_set[,9])/258
7. error_rf <- sum(predict(model_rf, newdata = test_set[, -9])!=
8.                 test_set[,9])/258
9. error_nn <- sum((predict(model_nn, newdata = test_set[, -9])[,2]>0.5)!=
10.                (test_set[,9]=="pos"))/258
11. error_ada <- sum((predict(model_ada, newdata = test_set[, -9])$prob[,2]>0.5)!=
12.                 (test_set[,9]=="pos"))/258
13. error_knn <- sum(fitted.values(model_knn)!=
14.                  test_set[,9])/258
15. error_gam <- sum((predict(model_gam, newdata = test_set[, -9])>1.5)!=
16.                 (test_set[,9]=="pos"))/258
17. print(c(error_tree,error_bagging,error_rf,error_nn,error_ada,error_knn,error_gam))
```

RESULT:

```
1. >print(c(error_tree,error_bagging,error_rf,error_nn,error_ada,error_knn,error_gam))
2. [1] 0.2558140 0.2480620 0.2325581 0.3449612 0.2596899 0.2790698 0.2403101
```

It seems that nearly all models have similar performance except **neural networks**, since the amount of data is too small to apply such an algorithm.

According to their performance and their interpretability, **I would like to choose classification tree** to analyse the data because of its **high interpretability**, although it does not have the best performance.

We are considering the reality.

Notice that the ability of predictions in such case is not so important, since it could pretty easy to check whether a person has diabetes. One could just go to hospital, and within one day, he/she could get the result of whether he/she has diabetes.

But instead, in reality, we are eager to figure out the relationship between the covariates and the response, which means that the interpretability are much more important than the prediction performance.

Then the choice could be pretty easy, since **classification tree** is much more interpretable than other models, especially when **bagging, random forest, neural network, Adaboost** and **knn** are nearly impossible to be interpreted and **GAM** are much harder to be interpreted than **classification tree**.

Task 5

By following the requirements of this task, we could easily change the data into that new dataset.

CODE:

```
1. data(PimaIndiansDiabetes2)
2. dat<-PimaIndiansDiabetes2[complete.cases(PimaIndiansDiabetes2),] #remove all NA obs
```

Then, we could just run the same code as before, and finally, the errors of each models could be:

```
1. >print(c(error_tree,error_bagging,error_rf,error_nn,error_ada,error_knn,error_gam))
2. [1] 0.2348485 0.1893939 0.1893939 0.3257576 0.2500000 0.2196970 0.2348485
```

COMMENTS:

Compared to the result of Task 4, we could easily find out that all errors are reduced with the new data.

This could be pretty understandable, since with changing the missing values into 0, we were actually creating wrong data and adding “noise” to the dataset. It could strongly influence the implementation of our model, since we are using biased data. But now, since we have no more “noise”, we are having a more precise result.

____END____