Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

CHAPTER-6

8.

(a)

```
1  set.seed(1)
2  X <- rnorm(100)
3  noise <- rnorm(100)
4
```

(b)

```
4
5  Y <- 3 + 1*X + 4*X^2 - 1*X^3 + noise
6
```

( c)

Regsubsets chooses three as the optimal number of parameters

```
12  require(leaps)
13  df <- data.frame(Y, X)
14  fit <- regsubsets(Y ~ poly(X, 10), data = df, nvmax = 10)
15
16  fit_summary <- summary(fit)
17
18  require(tidyverse)
19  require(ggplot2)
20  require(ggthemes)
21
22  data_frame(Cp = fit_summary$cp,
23            BIC = fit_summary$bic,
24            AdjR2 = fit_summary$adjr2) %>%
25    mutate(id = row_number()) %>%
26    gather(value_type, value, -id) %>%
27    ggplot(aes(id, value, col = value_type)) +
28    geom_line() + geom_point() + ylab('') + xlab('Number of Variables Used') +
29    facet_wrap(~ value_type, scales = 'free') +
30    theme_tufte() + scale_x_continuous(breaks = 1:10)
31
32
11:1   (Top Level)                                                    R Script
```

```
Console   Terminal   Jobs
~/
+     geom_line() + geom_point() + ylab('') + xlab('Number of Variables Used') +
+     facet_wrap(~ value_type, scales = 'free') +
+     theme_tufte() + scale_x_continuous(breaks = 1:10)
Warning message:
`data_frame()` was deprecated in tibble 1.1.0.
Please use `tibble()` instead.
This warning is displayed once every 8 hours.
Call `lifecycle::last_warnings()` to see where this warning was generated.
>
>
>
> |
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(d)
The backward and forward stepwise model agrees with the best subsets model

```
~/
> model <- train(Y ~ poly(X, 10), data = df,
+                method = 'glmStepAIC', direction = 'backward',
+                trace = 0,
+                trControl = trainControl(method = 'none', verboseIter = FALSE))
> postResample(predict(model, df), df$Y)
     RMSE   Rsquared        MAE
0.9314956 0.9569843 0.7488821
> summary(model$finalModel)

Call:
NULL

Deviance Residuals:
    Min      1Q  Median      3Q     Max
-1.8914 -0.5860 -0.1516  0.5892  2.1794

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.10265    0.09557  63.856  < 2e-16 ***
`poly(X, 10)1` -7.19295    0.95569  -7.526 2.96e-11 ***
`poly(X, 10)2` 40.74405    0.95569  42.633  < 2e-16 ***
`poly(X, 10)3` -14.70908   0.95569 -15.391  < 2e-16 ***
`poly(X, 10)5`  1.48019    0.95569   1.549    0.125
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 0.9133516)

    Null deviance: 2017.132  on 99  degrees of freedom
Residual deviance:   86.768  on 95  degrees of freedom
AIC: 281.59

Number of Fisher Scoring iterations: 2
```

```
> x_poly <- poly(df$X, 10)
> colnames(x_poly) <- paste0('poly', 1:10)
> model_forw <- train(y = Y, x = x_poly,
+                method = 'glmStepAIC', direction = 'forward',
+                trace = 0,
+                trControl = trainControl(method = 'none', verboseIter = FALSE))
> postResample(predict(model_forw, data.frame(x_poly)), df$Y)
     RMSE   Rsquared        MAE
0.9314956 0.9569843 0.7488821
> summary(model_forw$finalModel)

Call:
NULL

Deviance Residuals:
    Min      1Q  Median      3Q     Max
-1.8914 -0.5860 -0.1516  0.5892  2.1794

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    6.10265    0.09557  63.856  < 2e-16 ***
poly2         40.74405    0.95569  42.633  < 2e-16 ***
poly3        -14.70908    0.95569 -15.391  < 2e-16 ***
poly1         -7.19295    0.95569  -7.526 2.96e-11 ***
poly5          1.48019    0.95569   1.549    0.125
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 0.9133516)

    Null deviance: 2017.132  on 99  degrees of freedom
Residual deviance:   86.768  on 95  degrees of freedom
AIC: 281.59

Number of Fisher Scoring iterations: 2
```
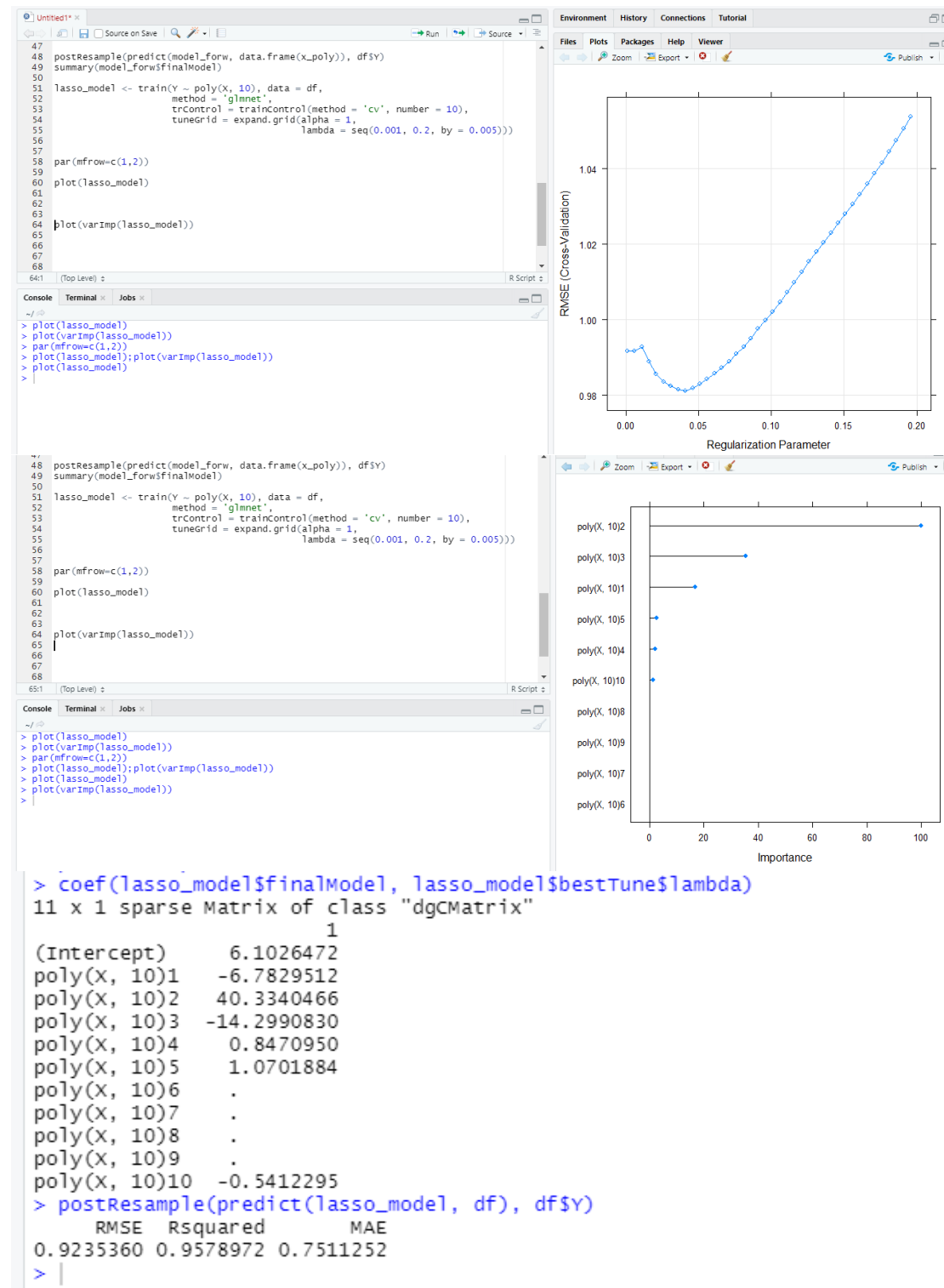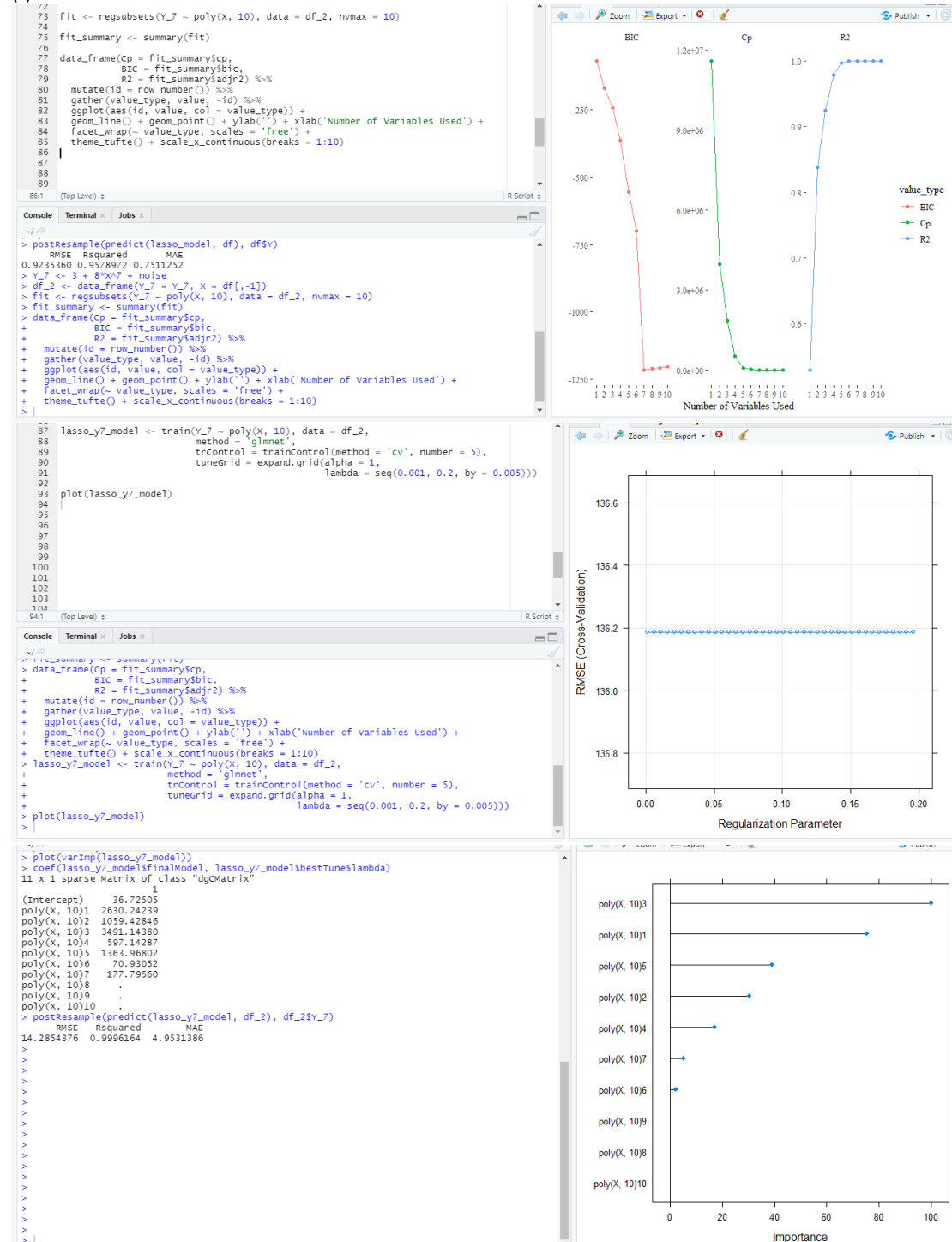
Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

( e)

The Lasso model overestimates the number of predictors needed. This might be expected since we used only RSS to select the optimal model but not the Bayesian Inference Criterion and Adjusted $R^2$ as regsubsets does or the Aikake Information Criterion as the stepwise selection does.





```
> coef(lasso_model$finalModel, lasso_model$bestTune$lambda)
11 x 1 sparse Matrix of class "dgCMatrix"
                        1
(Intercept)     6.1026472
poly(X, 10)1   -6.7829512
poly(X, 10)2   40.3340466
poly(X, 10)3  -14.2990830
poly(X, 10)4    0.8470950
poly(X, 10)5    1.0701884
poly(X, 10)6    .
poly(X, 10)7    .
poly(X, 10)8    .
poly(X, 10)9    .
poly(X, 10)10  -0.5412295
> postResample(predict(lasso_model, df), df$Y)
     RMSE  Rsquared       MAE
0.9235360 0.9578972 0.7511252
>
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(f)

```
73  fit <- regsubsets(Y_7 ~ poly(X, 10), data = df_2, nvmax = 10)
74
75  fit_summary <- summary(fit)
76
77  data_frame(Cp = fit_summary$cp,
78             BIC = fit_summary$bic,
79             R2 = fit_summary$adjr2) %>%
80    mutate(id = row_number()) %>%
81    gather(value_type, value, -id) %>%
82    ggplot(aes(id, value, col = value_type)) +
83    geom_line() + geom_point() + ylab('') + xlab('Number of Variables Used') +
84    facet_wrap(~ value_type, scales = 'free') +
85    theme_tufte() + scale_x_continuous(breaks = 1:10)
86  |
87
88
89
```
86:1   (Top Level) ÷                                                    R Script ÷

Console  Terminal ×  Jobs ×

```
> postResample(predict(lasso_model, df), df$Y)
      RMSE   Rsquared        MAE
 0.9235360  0.9578972  0.7511252
> Y_7 <- 3 + 8*X^7 + noise
> df_2 <- data_frame(Y_7 = Y_7, X = df[,-1])
> fit <- regsubsets(Y_7 ~ poly(X, 10), data = df_2, nvmax = 10)
> fit_summary <- summary(fit)
> data_frame(Cp = fit_summary$cp,
+            BIC = fit_summary$bic,
+            R2 = fit_summary$adjr2) %>%
+   mutate(id = row_number()) %>%
+   gather(value_type, value, -id) %>%
+   ggplot(aes(id, value, col = value_type)) +
+   geom_line() + geom_point() + ylab('') + xlab('Number of Variables Used') +
+   facet_wrap(~ value_type, scales = 'free') +
+   theme_tufte() + scale_x_continuous(breaks = 1:10)
> |
```



```
87  lasso_y7_model <- train(Y_7 ~ poly(X, 10), data = df_2,
88                          method = 'glmnet',
89                          trControl = trainControl(method = 'cv', number = 5),
90                          tuneGrid = expand.grid(alpha = 1,
91                                  lambda = seq(0.001, 0.2, by = 0.005)))
92
93  plot(lasso_y7_model)
94  |
95
96
97
98
99
100
101
102
103
104
```
94:1   (Top Level) ÷                                                    R Script ÷

Console  Terminal ×  Jobs ×

```
> fit_summary <- summary(fit)
> data_frame(Cp = fit_summary$cp,
+            BIC = fit_summary$bic,
+            R2 = fit_summary$adjr2) %>%
+   mutate(id = row_number()) %>%
+   gather(value_type, value, -id) %>%
+   ggplot(aes(id, value, col = value_type)) +
+   geom_line() + geom_point() + ylab('') + xlab('Number of Variables Used') +
+   facet_wrap(~ value_type, scales = 'free') +
+   theme_tufte() + scale_x_continuous(breaks = 1:10)
> lasso_y7_model <- train(Y_7 ~ poly(X, 10), data = df_2,
+                         method = 'glmnet',
+                         trControl = trainControl(method = 'cv', number = 5),
+                         tuneGrid = expand.grid(alpha = 1,
+                                 lambda = seq(0.001, 0.2, by = 0.005)))
> plot(lasso_y7_model)
> |
```



```
> plot(varImp(lasso_y7_model))
> coef(lasso_y7_model$finalModel, lasso_y7_model$bestTune$lambda)
11 x 1 sparse Matrix of class "dgCMatrix"
                       1
(Intercept)     36.72505
poly(X, 10)1  2630.24239
poly(X, 10)2  1059.42846
poly(X, 10)3  3491.14380
poly(X, 10)4   597.14287
poly(X, 10)5  1363.96802
poly(X, 10)6    70.93052
poly(X, 10)7   177.79560
poly(X, 10)8         .
poly(X, 10)9         .
poly(X, 10)10        .
> postResample(predict(lasso_y7_model, df_2), df_2$Y_7)
      RMSE   Rsquared        MAE
14.2854376  0.9996164  4.9531386
>
>
>
>
>
>
>
>
>
>
>
>
>
>
> |
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

9.
(a)
```r
library(ISLR)
library(caret)
library(tidyverse)
data('College')
set.seed(1)

inTrain <- createDataPartition(College$Apps, p = 0.75, list = FALSE)

training <- College[inTrain,]
testing <- College[-inTrain,]

Obj <- preProcess(training, method = c('center', 'scale'))

training <- predict(Obj, training)
testing <- predict(Obj, testing)

y_train <- training$Apps
y_test <- testing$Apps

apps <- dummyVars(Apps ~ ., data = training)
x_train <- predict(apps, training)
x_test <- predict(apps, testing)
```
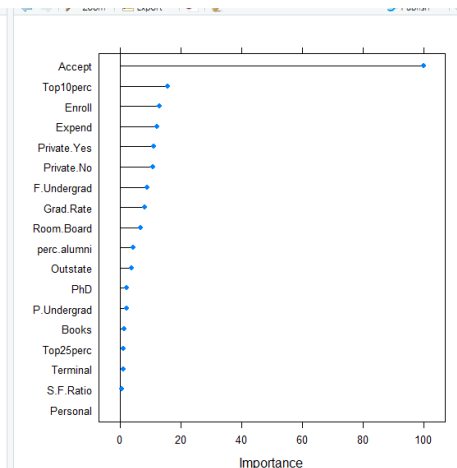
(b)
```r
model <- lm(Apps ~ ., data = training)
pred <- predict(model, testing)
info <- postResample(pred, testing$Apps)
```
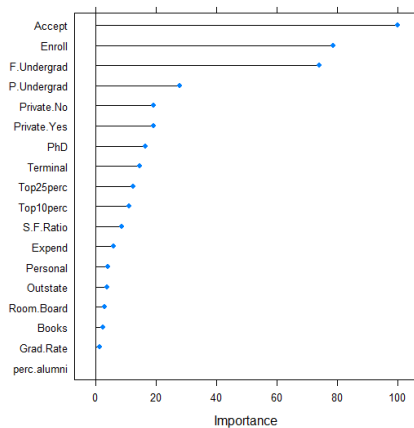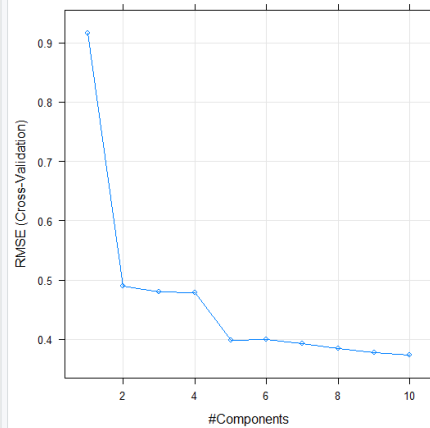
( c)
```r
> fit <- train(x = x_train, y = y_train,
+              method = 'glmnet',
+              trControl = trainControl(method = 'cv', number = 10),
+              tuneGrid = expand.grid(alpha = 0,
+                                     lambda = seq(0, 10e2, length.out = 20)))
warning message:
In nominalTrainWorkflow(x = x, y = y, wts = weights, info = trainInfo,  :
  There were missing values in resampled performance measures.
> (info <- postResample(predict(ridge_fit, x_test), y_test))
     RMSE  Rsquared       MAE
0.2853247 0.9211286 0.1645806
> coef(fit$finalModel, fit$bestTune$lambda)
19 x 1 sparse Matrix of class "dgCMatrix"
                      1
(Intercept)   0.034871314
Private.No    0.075423210
Private.Yes  -0.076037580
Accept        0.665628733
Enroll        0.090243372
Top10perc     0.107160248
Top25perc     0.011628030
F.Undergrad   0.063308801
P.Undergrad   0.017427317
Outstate     -0.028995432
Room.Board    0.048720533
Books         0.012799145
Personal     -0.002894430
PhD          -0.017989250
Terminal     -0.010434665
S.F.Ratio     0.006920126
perc.alumni  -0.031683867
Expend        0.083525070
Grad.Rate     0.058131023
> plot(ridge_fit)
> plot(varImp(ridge_fit))
>
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4
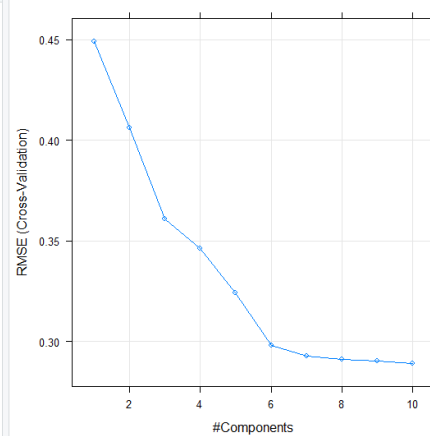
(d)

```
> fit <- train(x = x_train, y = y_train,
+                   method = 'glmnet',
+                   trControl = trainControl(method = 'cv', number = 10),
+                   tuneGrid = expand.grid(alpha = 1,
+                                          lambda = seq(0.0001, 1, length.out = 50)))
Warning message:
In nominalTrainWorkflow(x = x, y = y, wts = weights, info = trainInfo,  :
  There were missing values in resampled performance measures.
> (info <- postResample(predict(fit, x_test), y_test))
     RMSE   Rsquared        MAE
0.2802352 0.9201823 0.1561301
> coef(fit$finalModel, fit$bestTune$lambda)
19 x 1 sparse Matrix of class "dgCMatrix"
                      1
(Intercept) -0.037243095
Private.No   0.137026483
Private.Yes  .
Accept       1.041851224
Enroll      -0.202744295
Top10perc    0.201576111
Top25perc   -0.046294002
F.Undergrad  0.012507811
P.Undergrad  0.029491934
Outstate    -0.085333127
Room.Board   0.033826427
Books        0.005116779
Personal     0.006295093
PhD         -0.037015361
Terminal    -0.002461461
S.F.Ratio    0.005385825
perc.alumni -0.006575661
Expend       0.077037030
Grad.Rate    0.037985756
> plot(fit)
>
```

```
+                   trControl = trainControl(method = 'cv', number = 10),
+                   tuneGrid = expand.grid(alpha = 1,
+                                          lambda = seq(0.0001, 1, length.out = 50)))
Warning message:
In nominalTrainWorkflow(x = x, y = y, wts = weights, info = trainInfo,  :
  There were missing values in resampled performance measures.
> (info <- postResample(predict(fit, x_test), y_test))
     RMSE   Rsquared        MAE
0.2802352 0.9201823 0.1561301
> coef(fit$finalModel, fit$bestTune$lambda)
19 x 1 sparse Matrix of class "dgCMatrix"
                      1
(Intercept) -0.037243095
Private.No   0.137026483
Private.Yes  .
Accept       1.041851224
Enroll      -0.202744295
Top10perc    0.201576111
Top25perc   -0.046294002
F.Undergrad  0.012507811
P.Undergrad  0.029491934
Outstate    -0.085333127
Room.Board   0.033826427
Books        0.005116779
Personal     0.006295093
PhD         -0.037015361
Terminal    -0.002461461
S.F.Ratio    0.005385825
perc.alumni -0.006575661
Expend       0.077037030
Grad.Rate    0.037985756
> plot(fit)
> plot(varImp(fit))
>
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

( e)

```
package 'pls' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
        C:\Users\Santhoshini sree\AppData\Local\Temp\Rtmp4e4wQZ\downloaded_packages
> pcr_model <- train(x = x_train, y = y_train,
+                    method = 'pcr',
+                    trControl = trainControl(method = 'cv', number = 10),
+                    tuneGrid = expand.grid(ncomp = 1:10))
> (pcr_info <- postResample(predict(pcr_model, x_test), y_test))
    RMSE  Rsquared       MAE
0.3231292 0.8916531 0.1986075
> coef(pcr_model$finalModel)
, , 10 comps

               .outcome
Private.No    0.031985972
Private.Yes  -0.031985972
Accept        0.343576750
Enroll        0.305359773
Top10perc     0.042630417
Top25perc     0.027790893
F.Undergrad   0.273818439
P.Undergrad  -0.049487667
Outstate      0.038573119
Room.Board    0.070607615
Books         0.016433593
Personal     -0.023529455
PhD          -0.023992433
Terminal     -0.024182230
S.F.Ratio     0.003741623
perc.alumni  -0.070567887
Expend        0.090126298
Grad.Rate     0.071302714

> plot(pcr_model)
>
173:1   (Top Level)                                                        R Script
```

```
Console   Terminal   Jobs
~/
> pcr_model <- train(x = x_train, y = y_train,
+                    method = 'pcr',
+                    trControl = trainControl(method = 'cv', number = 10),
+                    tuneGrid = expand.grid(ncomp = 1:10))
> (pcr_info <- postResample(predict(pcr_model, x_test), y_test))
    RMSE  Rsquared       MAE
0.3231292 0.8916531 0.1986075
> coef(pcr_model$finalModel)
, , 10 comps

               .outcome
Private.No    0.031985972
Private.Yes  -0.031985972
Accept        0.343576750
Enroll        0.305359773
Top10perc     0.042630417
Top25perc     0.027790893
F.Undergrad   0.273818439
P.Undergrad  -0.049487667
Outstate      0.038573119
Room.Board    0.070607615
Books         0.016433593
Personal     -0.023529455
PhD          -0.023992433
Terminal     -0.024182230
S.F.Ratio     0.003741623
perc.alumni  -0.070567887
Expend        0.090126298
Grad.Rate     0.071302714

> plot(pcr_model)
> plot(varImp(pcr_model))
>
```

(f)

```
~/
> model <- train(x = x_train, y = y_train,
+                method = 'pls',
+                trControl = trainControl(method = 'cv', number = 10),
+                tuneGrid = expand.grid(ncomp = 1:10))
> (info <- postResample(predict(model, x_test), y_test))
    RMSE  Rsquared       MAE
0.2792580 0.9209302 0.1572165
> coef(model$finalModel)
, , 10 comps

               .outcome
Private.No    0.071314109
Private.Yes  -0.071314109
Accept        1.039263053
Enroll       -0.169855531
Top10perc     0.235572152
Top25perc    -0.070270182
F.Undergrad  -0.022580399
P.Undergrad   0.032522451
Outstate     -0.085919942
Room.Board    0.036344367
Books         0.004835496
Personal      0.007546695
PhD          -0.044709439
Terminal      0.002332025
S.F.Ratio     0.010868271
perc.alumni  -0.009046689
Expend        0.072963738
Grad.Rate     0.037425850

> plot(model)
>
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(g)

```
> as_data_frame(rbind(lin_info,
+                     ridge_info,
+                     lasso_info,
+                     pcr_info,
+                     pls_info)) %>%
+   mutate(model = c('Linear', 'Ridge', 'Lasso', 'PCR', 'PLS')) %>%
+   select(model, RMSE, Rsquared)
# A tibble: 5 x 3
  model    RMSE Rsquared
  <chr>   <dbl>    <dbl>
1 Linear  0.280    0.920
2 Ridge   0.285    0.921
3 Lasso   0.291    0.914
4 PCR     0.323    0.892
5 PLS     0.279    0.921
```

```
3 Lasso  0.291  0.914
4 PCR    0.323  0.892
5 PLS    0.279  0.921
Warning message:
`as_data_frame()` was deprecated in tibble 2.0.0.
Please use `as_tibble()` instead.
The signature and semantics have changed, see `?as_tibble`.
This warning is displayed once every 8 hours.
Call `lifecycle::last_warnings()` to see where this warning was generated.
> testing %>%
+   summarize(sd = sd(Apps))
        sd
1 0.9818241
> library(ggthemes)
> residfunc <- function(fit, data) {
+   predict(fit, data) - testing$Apps
+ }
> data_frame(observed = testing$Apps,
+            LM = residfunc(lin_model, testing),
+            Ridge = residfunc(ridge_fit, x_test),
+            Lasso = residfunc(lasso_fit, x_test),
+            PCR = residfunc(pcr_model, x_test),
+            PLS = residfunc(pls_model, x_test)) %>%
+   gather(Model, Residuals, -Observed) %>%
+   ggplot(aes(observed, Residuals, col = Model)) +
+   geom_hline(yintercept = 0, lty = 2) +
+   geom_point(alpha = 0.6) +
+   geom_smooth(method = 'loess', alpha = 0.01, col = 'lightsalmon2') +
+   facet_wrap(~ Model, ncol = 5) +
+   theme_tufte() +
+   theme(legend.position = 'top') +
+   coord_flip()
`geom_smooth()` using formula 'y ~ x'
Warning message:
attributes are not identical across measure variables;
they will be dropped
> |
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

10.
(a)
```
> set.seed(1)
> x <- matrix(rnorm(1000 * 20), 1000, 20)
> b <- rnorm(20)
> b[3] <- 0
> b[4] <- 0
> b[9] <- 0
> b[19] <- 0
> b[10] <- 0
> eps <- rnorm(1000)
> y <- x %*% b + eps
>
```

(b)

Console   Terminal ×   Jobs ×

~/ 
```
> train <- sample(seq(1000), 100, replace = FALSE)
> test <- -train
> x.train <- x[train, ]
> x.test <- x[test, ]
> y.train <- y[train]
> y.test <- y[test]
>
```

( c)

Console   Terminal ×   Jobs ×

~/ 
```
> data.train <- data.frame(y = y.train, x = x.train)
> library(leaps)
> regfit.full <- regsubsets(y ~ ., data = data.train, nvmax = 20)
> train.mat <- model.matrix(y ~ ., data = data.train, nvmax = 20)
> val.errors <- rep(NA, 20)
> for (i in 1:20) {
+    coefi <- coef(regfit.full, id = i)
+    pred <- train.mat[, names(coefi)] %*% coefi
+    val.errors[i] <- mean((pred - y.train)^2)
+ }
>
```
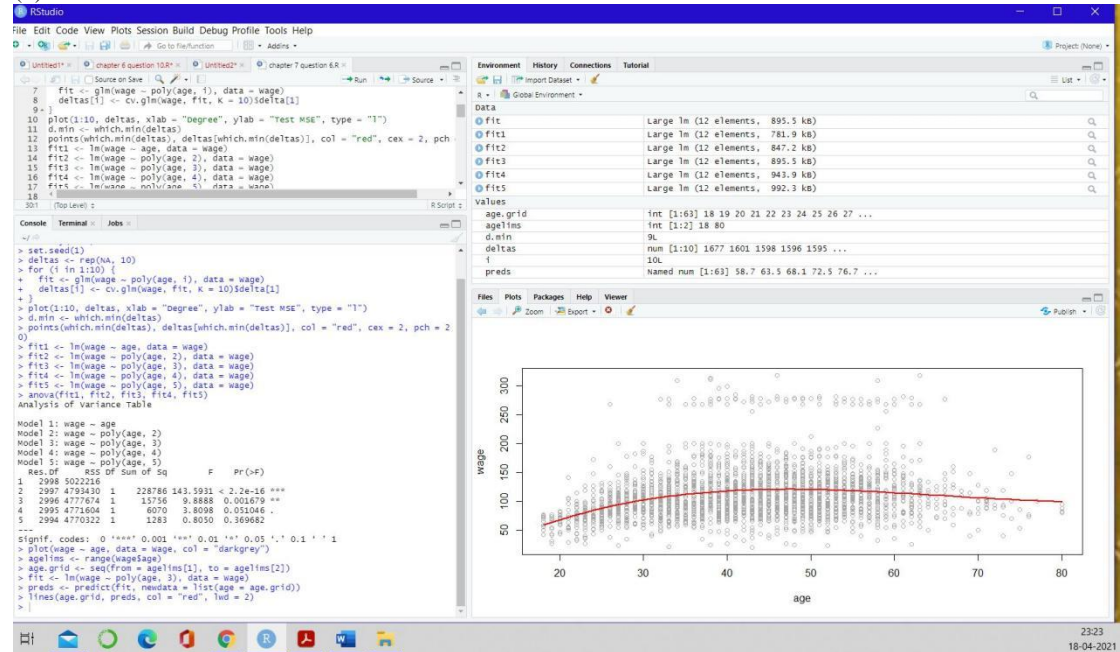
Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

```
> regfit.full <- regsubsets(y ~ ., data = data.train, nvmax = 20)
> train.mat <- model.matrix(y ~ ., data = data.train, nvmax = 20)
> val.errors <- rep(NA, 20)
> for (i in 1:20) {
+    coefi <- coef(regfit.full, id = i)
+    pred <- train.mat[, names(coefi)] %*% coefi
+    val.errors[i] <- mean((pred - y.train)^2)
+ }
> plot(val.errors, xlab = "Number of predictors", ylab = "Training MSE", pch = 19, type = "b")
>
```



(d)

```
> data.test <- data.frame(y = y.test, x = x.test)
> test.mat <- model.matrix(y ~ ., data = data.test, nvmax = 20)
> val.errors <- rep(NA, 20)
> for (i in 1:20) {
+    coefi <- coef(regfit.full, id = i)
+    pred <- test.mat[, names(coefi)] %*% coefi
+    val.errors[i] <- mean((pred - y.test)^2)
+ }
> plot(val.errors, xlab = "Number of predictors", ylab = "Test MSE", pch = 19, type = "b")
>
```



(d)

```
> data.test <- data.frame(y = y.test, x = x.test)
> test.mat <- model.matrix(y ~ ., data = data.test, nvmax = 20)
> val.errors <- rep(NA, 20)
> for (i in 1:20) {
+    coefi <- coef(regfit.full, id = i)
+    pred <- test.mat[, names(coefi)] %*% coefi
+    val.errors[i] <- mean((pred - y.test)^2)
+ }
> plot(val.errors, xlab = "Number of predictors", ylab = "Test MSE", pch = 19, type = "b")
>
```

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(e )

```
> which.min(val.errors)
[1] 19
>
```

14-variables model has the smallest test MSE.

(f)

```
> coef(regfit.full, which.min(val.errors))
 (Intercept)          x.1          x.2          x.3          x.5          x.6          x.7
 -0.04011778   0.13698928   0.20775097  -0.10438563   1.04174837  -0.21709681  -1.31761521
         x.8          x.9         x.10         x.11         x.12         x.13         x.14
  0.72571564   0.12387155  -0.18363204   1.01888399   0.64149490  -0.41714902  -0.70005302
        x.15         x.16         x.17         x.18         x.19         x.20
 -0.76105664  -0.40810077   0.04686190   1.65688296  -0.13786948  -0.99777611
>
```

The best model caught all zeroed out coefficients.

(g)
We may see that the model with 3 variables minimizes the error between the estimated and true coefficients. However test error is minimized by the model with 14 variables. So, a better fit of true coefficients doesn't necessarily mean a lower test MSE.

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

CHAPTER 7

6.

(a)

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(b)

Alexa Summers
Santhoshini Sree Bolisetty
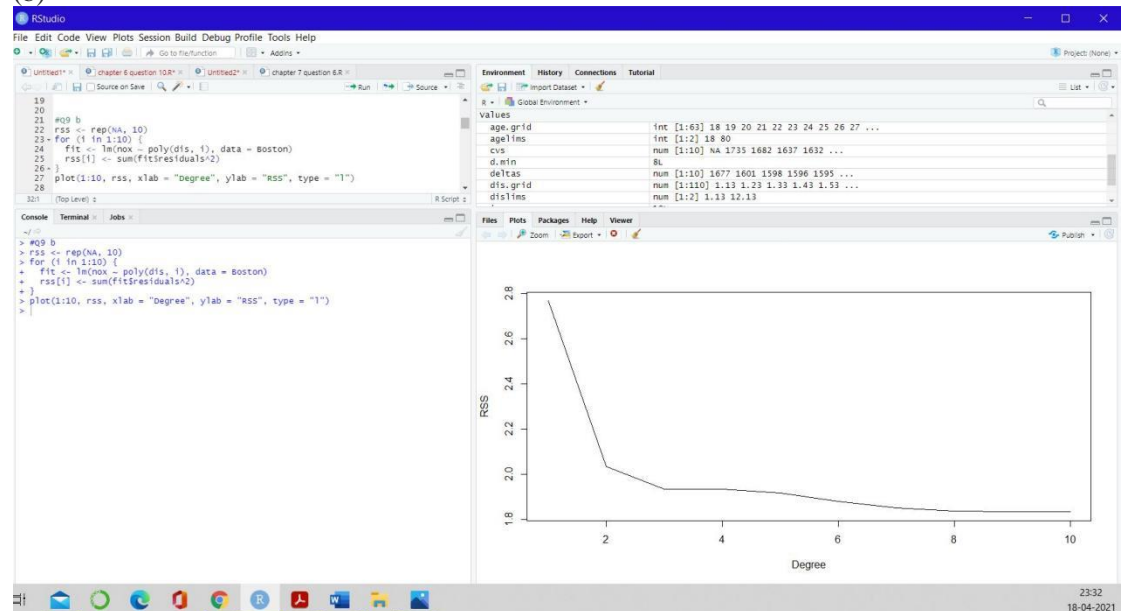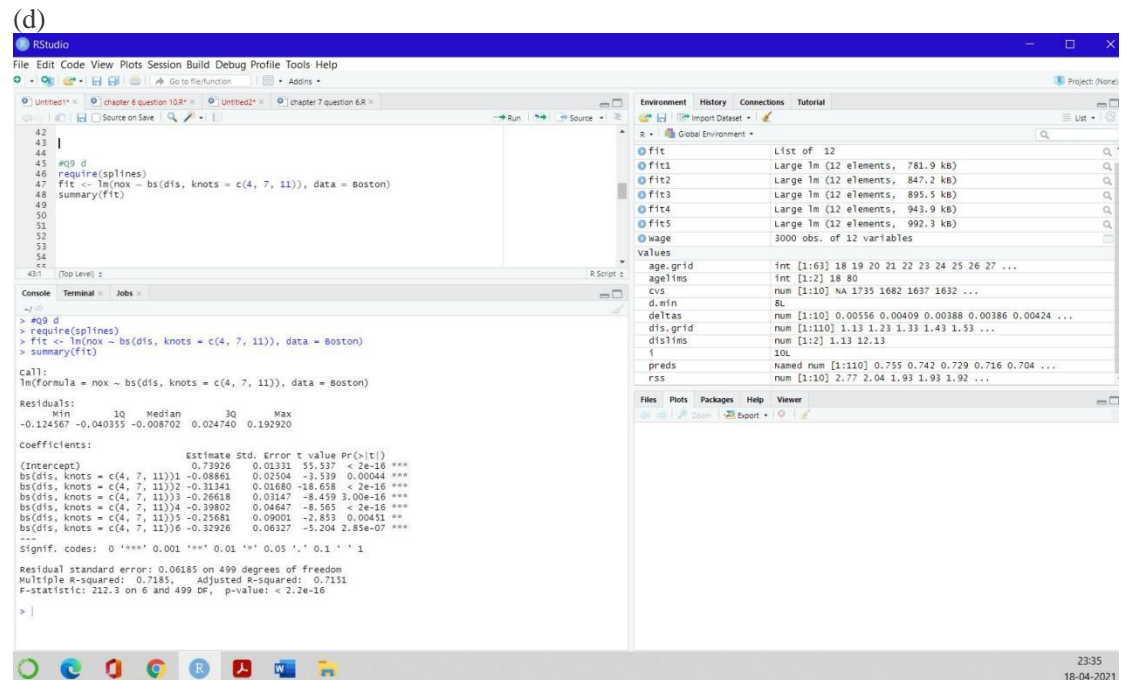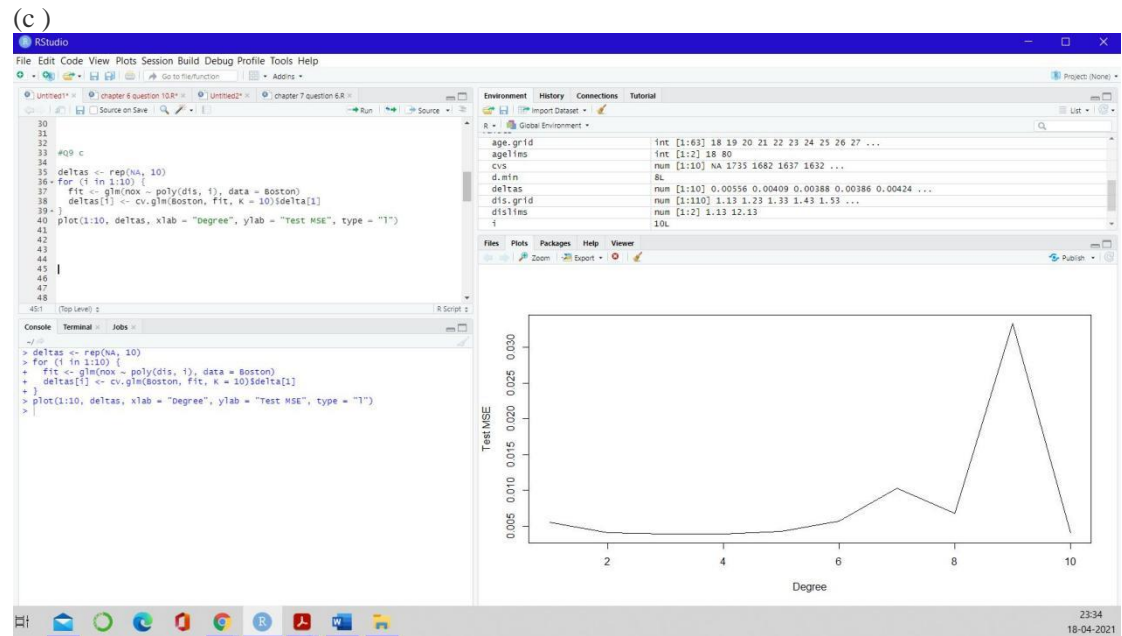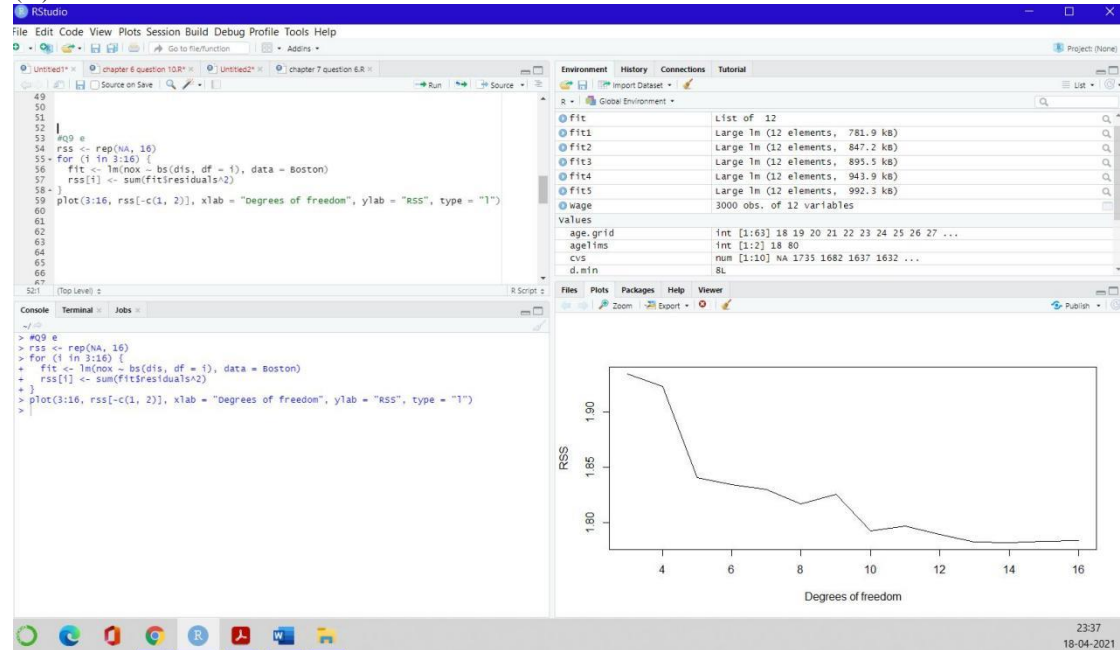Gireesh Kumar Muppalla
CS 5565—Lab 4

9.

(a)



(b)

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

(c)



(d)

Alexa Summers
Santhoshini Sree Bolisetty
Gireesh Kumar Muppalla
CS 5565—Lab 4

( e)



(f)