

PA_HW5

Group10

2022-11-08

Problem 3

3(a)

```
library(MASS)
# data set
Boston$logcrim = log(Boston$crim) # create log transform of crim
# summary(Boston)
set.seed(12345)
train=runif(nrow(Boston))<= .5 # pick train/test split
print("Frequency Distribution")

## "Frequency Distribution"

table(train) # frequency distribution

## train
## FALSE TRUE
## 282 224
```

3(b)

```
fit1 = lm(logcrim ~ . - crim, data = Boston, subset=train)
summary(fit1)
p1 = predict(fit1, Boston[!train,])
# MSE
mse_test = mean((Boston$logcrim[train] - p1)^2)
mse_test

## [1] 0.7083435

# residual plot
plot(fit1, which = 1, pch = 16, cex = .5)
```



the residual plot we learn that the variance are not very constant.

3(c)

```
fit2 = step(fit1, trace = F)
p2 = predict(fit2, Boston[!train,])
mean((Boston$logcrim[train] - p2)^2)

## [1] 0.7033381

The test set MSE is 0.7033381.
```

3(d)

```
library("glmnet")

## Warning: package 'glmnet' was built under R version 4.2.2

## Loading required package: Matrix

## Loaded glmnet 4.1-4

# fit ridge
x = model.matrix(logcrim ~ . - crim -1, Boston)
fit_ridge = glmnet(x[train,], Boston$logcrim[train], alpha = 0)
# CV
set.seed(1234)
ridge_cv = cv.glmnet(x[train,], Boston$logcrim[train], alpha = 0) # find optimal lambda
l1 = ridge_cv$lambda.min
l1

## [1] 0.1866301

# test set mse
p3 = predict(fit_ridge, s=l1, newx=x[!train,])
mean((Boston$logcrim[train] - p3)^2)

## [1] 0.7767352

The test set MSE is 0.7767352.
```

3(e)

```
# fit lasso
x = model.matrix(logcrim ~ . - crim -1, Boston)
fit_lasso = glmnet(x[train,], Boston$logcrim[train], alpha = 1)
# CV
set.seed(1234)
lasso_cv = cv.glmnet(x[train,], Boston$logcrim[train], alpha = 1) # find optimal lambda
l2 = lasso_cv$lambda.min
l2

## [1] 0.009288624

# test set mse
p4 = predict(fit_lasso, s=l2, newx=x[!train,])
mean((Boston$logcrim[train] - p4)^2)

## [1] 0.7023993

The test set MSE is 0.7023993.
```

3(f)

```
# ridge model
x = model.matrix(logcrim ~zn
+indus+I(indus^2)
+chas
+nox+I(nox^2)
+rm+log(rm)
+age+log(age)
+dis+log(dis)
+rad+log(rad)
+tax+log(tax)
+ptratio+log(ptratio)
+black+log(black)
+lstat+log(lstat)
+medv+log(medv), Boston)
fit_r = glmnet(x[train,], Boston$logcrim[train], alpha=0)
fit_cv = cv.glmnet(x[train,], Boston$logcrim[train], alpha=0)
l_r = fit_cv$lambda.min

yhat = predict(fit_r, s=fit_cv$lambda.min, newx=x[!train,])
mean((Boston$logcrim[train] - yhat)^2)

## [1] 0.6510401

# lasso model
fit_l1 = glmnet(x[train,], Boston$logcrim[train], alpha=1)
fit_cv2 = cv.glmnet(x[train,], Boston$logcrim[train], alpha=1)
l_l = fit_cv2$lambda.min

yhat2 = predict(fit_l1, s=fit_cv2$lambda.min, newx=x[!train,])
mean((Boston$logcrim[train] - yhat2)^2)

## [1] 0.5612213

# step-wise
fit4 = lm(logcrim ~ 1, Boston, subset=train)
fit_s = step(fit4, scopes=zn
+indus+I(indus^2)
+chas
+nox+I(nox^2)
+rm+log(rm)
+age+log(age)
+dis+log(dis)
+rad+log(rad)
+tax+log(tax)
+ptratio+log(ptratio)
+black+log(black)
+lstat+log(lstat)
+medv+log(medv), trace = FALSE)
yhat3 = predict(fit_s, Boston[!train,])
mean((Boston$logcrim[train] - yhat3)^2)

## [1] 0.5694736
```

Log transformations and square transformations of predictors are important. Because these transformations help reduce MSE in all three models.

Problem 4

4(a)

```
sigma = chol(0.9, nrow=4, ncol=4) + .1*diag(4)
A = chol(sigma)
t(A) %*% A

##      [,1] [,2] [,3] [,4]
## [1,] 1.0 0.0 0.0 0.0
## [2,] 0.0 1.0 0.0 0.0
## [3,] 0.0 0.0 1.0 0.0
## [4,] 0.0 0.0 0.0 1.0
```

4(b)

```
X = matrix(rnorm(4000), nrow=1000)
Z = Z %*% A
var(X)

##      [,1] [,2] [,3] [,4]
## [1,] 0.9219924 0.8492134 0.8528997 0.8503184
## [2,] 0.8492134 0.7769810 0.8667100 0.8601371
## [3,] 0.8528997 0.8667100 0.9625422 0.8652925
## [4,] 0.8503184 0.8601371 0.8652925 0.9634800
```

It approximately equal \$\$.

4(c)

```
set.seed(12345)
# generate a new Z, A and X
n = 10100
ntest=n-ntrain
ntest=n-ntrain
Z <- matrix(rnorm(n*15), nrow=n)
sigma <- matrix(0.9, nrow=15, ncol=15) + diag(rep(1-0.9, 15))
A <- cho(sigm)
X <- Z %*% A

beta = c(1,-1,1,5,0.5,-0.5,rep(0,10))
e = rnorm(10100)*3
y = 3 + X %*% beta + e

train <- data.frame(X[1:ntrain,], y=y[1:ntrain])
test <- data.frame(X[ntrain+1:n,], y=y[ntrain+1:n])

mean((test$y-predict(fit, test))^2)

## [1] 9.449501

confint(fit)

##      2.5 %      97.5 %
## (Intercept) 2.3714294 3.6797538
## X1          -0.8182307 2.7059553
## X2          -3.6208313 0.3695036
## X3           1.0609970 4.5597365
## X4          -2.3769925 1.7692921
## X5          -1.9919717 1.2488126

## Residual standard error: 3.2 on 94 degrees of freedom
## Multiple R-squared: 0.2407, Adjusted R-squared: 0.2003
## F-statistic: 5.961 on 5 and 94 DF, p-value: 7.843e-05
```

The residual standard error is 3.2. Slopes for $x_1(-1.2, 3.4, 5)$ are 0.9439, -1.6256, 2.7079, -0.3034, -0.3711. R-squared is 0.2407. Slope for x_1 roughly equals the true parameter and the other slopes aren't. Slope for x_4 has the wrong sign, other slopes have the right signs. In terms of statistical significance, only slope for x_3 is significant at 0.01 level. All the 95% confidence interval covers the true value.

4(e)

```
mean((test$y-predict(fit, test))^2)

## [1] 9.449501

The value of MSE is 9.449501.
```

4(f)

```
fit_4f = lm(y ~ ., data = train)
summary_4f = summary(fit_4f)
coeffs = summary_4f$coefficients
betas_4f = coeffs[2:16]
se = coeffs[18:32]
ci_l = c()
ci_u = c()
for (idx in 1:length(betas_4f)) {
  ci_l = append(ci_l, (betas_4f[idx] - 2 * se[idx]))
  ci_u = append(ci_u, (betas_4f[idx] + 2 * se[idx]))
}

is_in = c()
for (idx in 1:length(betas_4f)) {
  if ((betas[idx] > ci_l[idx]) & (betas[idx] < ci_u[idx])) {
    is_in = append(is_in, 1)
  } else {
    is_in = append(is_in, 0)
  }
}
vars_in = sum(is_in)
print(summary(fit_4f))

##      Call:
## lm(formula = y ~ ., data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.7768 -2.0442  0.2997  1.8333  6.9526
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.0296      0.3295  9.183 1e-14 ***
## X1          0.9439      0.0075  1.064 0.29026
## X2         -1.6256      0.0075 -1.618 0.10906
## X3          2.7079      0.8924  3.124 0.00237 **
## X4         -0.3034      1.0439 -0.291 0.77280
## X5         -0.3711      0.8164 -0.455 0.65048
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 3.2 on 94 degrees of freedom
## Multiple R-squared: 0.2407, Adjusted R-squared: 0.2003
## F-statistic: 5.961 on 5 and 94 DF, p-value: 0.007036
```

- 15 of the 15 are in the confidence intervals
- The only feature that's significant is x_3 or β_3 . None of the other x_i (1 through 5) are significant.
- All of the signs are not correct for beta 1 through 5, β_5 should be positive.

4(g)

```
mse_4g = mean((test$y-predict(fit_4f, test))^2)
```

The MSE is 10.2347653 for the full model.

4(h)

```
step_4h = stepAIC(lm(y ~., data = train),
direction = "both",
trace = F)
summary(step_4h)

##
## Call:
## lm(formula = y ~ X1 + X3 + X8 + X13, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.2310 -1.8975  0.2254  1.6861  7.4489
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.09711      0.34378  9.009 5.69e-14 ***
## X1          1.64539      1.01410  1.622 0.10845
## X2         -1.27455      1.20322 -1.132 0.26102
## X3          3.04446      0.99629  3.056 0.00301 **
## X4          0.17894      1.16865  0.153 0.87867
## X5          0.12057      0.95410  0.126 0.89974
## X6         -0.42267      1.04928 -0.402 0.68808
## X7         -0.05058      1.04946 -0.048 0.96547
## X8         -1.48874      1.18517 -1.256 0.22255
## X9          1.02701      1.03928  0.988 0.32589
## X10         -0.83861      1.13590 -0.738 0.46179
## X11         0.68516      1.02708  0.667 0.50691
## X12         -0.55163      1.07908 -0.511 0.61055
## X13         0.52500      1.22391 -0.026 0.30773
## X14         0.25319      1.01348  0.250 0.80705
## X15         0.73817      1.23259 -0.599 0.55185
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 3.258 on 84 degrees of freedom
## Multiple R-squared: 0.2964, Adjusted R-squared: 0.1708
## F-statistic: 2.359 on 15 and 84 DF, p-value: 0.007036
```

- 15 of the 15 are in the confidence intervals
- The only feature that's significant is x_3 or β_3 . None of the other x_i (1 through 5) are significant.
- All of the signs are not correct for beta 1 through 5, β_5 should be positive.

4(i)

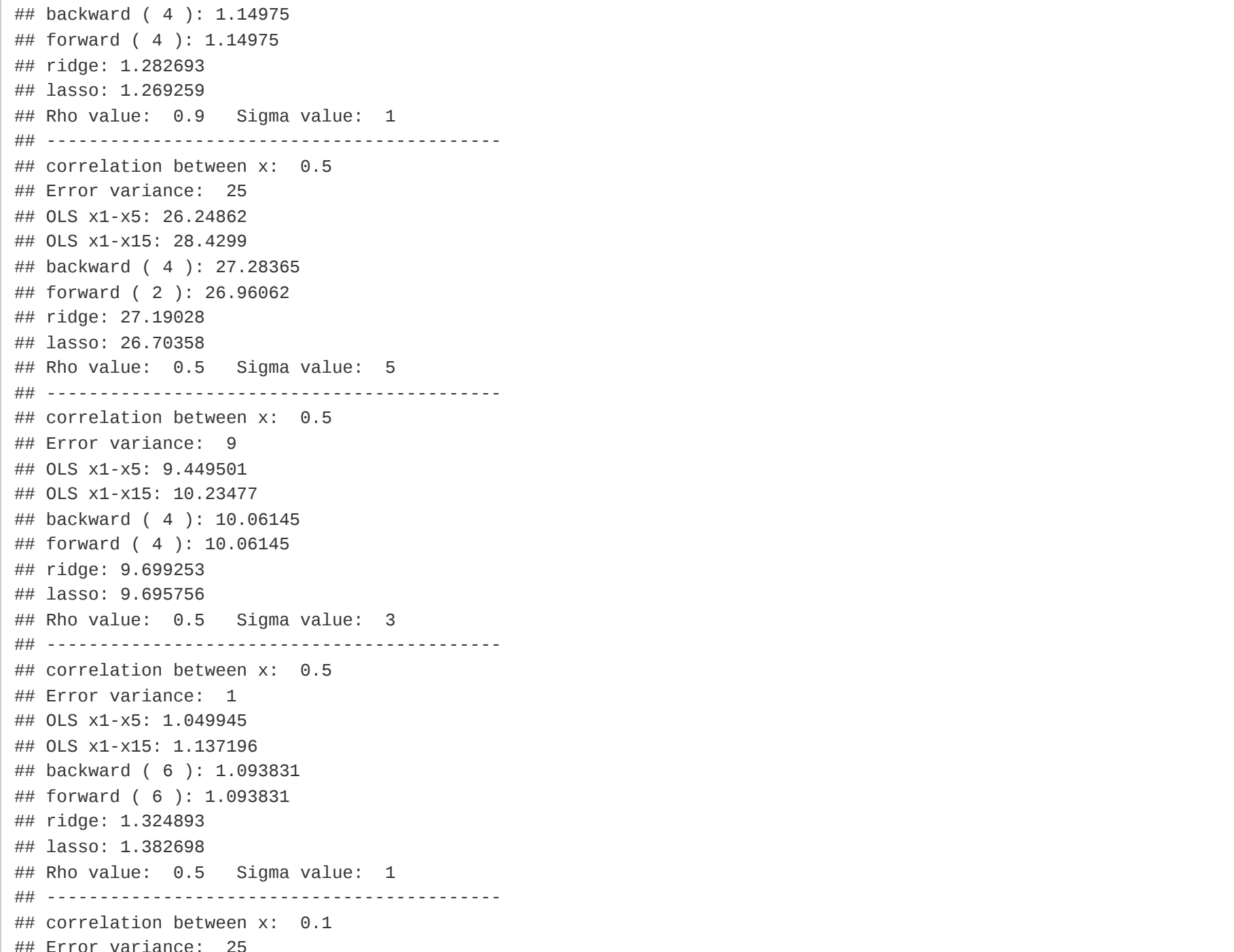
```
mse_4i = mean((test$y-predict(step_4h, test))^2)
```

The MSE is 10.0442626 for the model found using stepAIC.

4(j)

```
x_4j = model.matrix(y ~ . -1, data = train)
set.seed(1234)
lasso_cv_4j = cv.glmnet(x_4j, train$y, alpha = 0, nfolds = 5)
plot(lasso_cv_4j, xvar = "lambda", lwd = 1)

## Warning in plot.window(...): "xvar" is not a graphical parameter
##
## Warning in plot.window(...): "label" is not a graphical parameter
##
## Warning in plot.xy(xy, type, ...): "xvar" is not a graphical parameter
##
## Warning in plot.xy(xy, type, ...): "label" is not a graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "xvar" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "label" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "xvar" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "label" is not a
## graphical parameter
##
## Warning in box(...): "xvar" is not a graphical parameter
##
## Warning in box(...): "label" is not a graphical parameter
##
## Warning in title(...): "xvar" is not a graphical parameter
##
## Warning in title(...): "label" is not a graphical parameter
title("Ridge Trace")
```



4(k)

```
x_4k = model.matrix(y ~ . -1, data = test)
mse_4k = mean((test$y - predict(ridge_cv_4j, s = ridge_cv_4j$lambda.min, x_4k))^2)
```

The MSE is 9.3658854 for the Ridge model.

4(l)

```
x_4l = model.matrix(y ~ . -1, data = train)
set.seed(1234)
lasso_cv_4l = cv.glmnet(x_4l, train$y, alpha = 1, nfolds = 5)
plot(lasso_cv_4l, xvar = "lambda", lwd = 1)

## Warning in plot.window(...): "xvar" is not a graphical parameter
##
## Warning in plot.window(...): "label" is not a graphical parameter
##
## Warning in plot.xy(xy, type, ...): "xvar" is not a graphical parameter
##
## Warning in plot.xy(xy, type, ...): "label" is not a graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "xvar" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "label" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "xvar" is not a
## graphical parameter
##
## Warning in axis(side = side, at = at, labels = labels, ...): "label" is not a
## graphical parameter
##
## Warning in box(...): "xvar" is not a graphical parameter
##
## Warning in box(...): "label" is not a graphical parameter
##
## Warning in title(...): "xvar" is not a graphical parameter
##
## Warning in title(...): "label" is not a graphical parameter
title("Lasso Trace")
```


4(m)

```
x_4m = model.matrix(y ~ . -1, data = test)
mse_4m = mean((test$y - predict(lasso_cv_4l, s = lasso_cv_4l$lambda.min, x_4m))^2)
```

The MSE is 9.4173197 for the Ridge model.

4(n)

```
source("hw5.R")
rho = c(0.9,0.5,0.1)
sigma = c(5,3,1)

for (r in rho) {
  for (s in sigma) {
    hws(rho = r, sigmae = s)
    cat("Rho value: ", r, "\n", "Sigma value: ", s, "\n")
  }
}

## -----
## correlation between x: 0.9
## Error variance: 25
## OLS x1-x5: 26.24862
## OLS x1-x15: 28.4299
## backward ( 2 ): 26.96192
## forward ( 2 ): 26.96192
## lasso: 25.89624
## Rho value: 0.9 Sigma value: 5
## -----
## correlation between x: 0.9
## Error variance: 9
## OLS x1-x5: 9.449501
## OLS x1-x15: 10.23477
## backward ( 4 ): 10.04426
## forward ( 2 ): 9.879715
## ridge: 9.365885
## lasso: 9.41732
## Rho value: 0.9 Sigma value: 3
## -----
## correlation between x: 0.9
## Error variance: 25
## OLS x1-x5: 26.24862
## OLS x1-x15: 28.4299
## backward ( 4 ): 1.134975
## forward ( 4 ): 1.134975
## ridge: 1.282698
## lasso: 1.269259
## Rho value: 0.9 Sigma value: 1
## -----
## correlation between x: 0.5
## Error variance: 25
## OLS x1-x5: 26.24862
## OLS x1-x15: 28.4299
## backward ( 4 ): 27.23365
## forward ( 2 ): 26.96062
## ridge: 27.19028
## lasso: 26.70358
## Rho value: 0.5 Sigma value: 5
## -----
## correlation between x: 0.5
## Error variance: 9
## OLS x1-x5: 9.449501
## OLS x1-x15: 10.23477
## backward ( 4 ): 10.06145
## forward ( 4 ): 10.06145
## ridge: 9.498251
## lasso: 9.695756
## Rho value: 0.5 Sigma value: 3
## -----
## correlation between x: 0.5
## Error variance: 25
## OLS x1-x5: 26.24862
## OLS x1-x15: 28.4299
## backward ( 4 ): 1.093831
## forward ( 6 ): 1.093831
## ridge: 1.324893
## lasso: 1.382698
## Rho value: 0.5 Sigma value: 1
## -----
## correlation between x: 0.1
## Error variance: 25
## OLS x1-x5: 26.24862
## OLS x1-x15: 28.4299
## backward ( 4 ): 27.49461
## forward ( 2 ): 27.49461
## ridge: 27.6429
## lasso: 27.17986
## Rho value: 0.1 Sigma value: 5
## -----
## correlation between x: 0.1
## Error variance: 9
## OLS x1-x5: 9.449501
## OLS x1-x15: 10.23477
## backward ( 4 ): 10.27753
## forward ( 4 ): 10.27753
## ridge: 9.917381
## lasso: 9.858322
## Rho value: 0.1 Sigma value: 3
## -----
## correlation between x: 0.1
## Error variance: 1
## OLS x1-x5: 1.049945
## OLS x1-x15: 1.137196
## backward ( 4 ): 1.137196
## forward ( 6 ): 1.095551
## ridge: 1.31475
## lasso: 1.264195
## Rho value: 0.1 Sigma value: 1
```

When there is moderate to low rho/multicollinearity, OLS generally performs better than ridge or lasso regression. Ridge or lasso regression work better when there is high multicollinearity and moderate to high noise. When using lasso or ridge we will use lasso when noise is higher. In this test they perform best when rho is 0.9 and sigma is greater than or equal to 5.