

P8106 HW1

Shihui Zhu

Contents

| | |
|---|-----------|
| (a) Least Square | 1 |
| Perform 10-folds cross-validation on the full model | 2 |
| MSE, R-squared values | 3 |
| (b) Lasso | 3 |
| Fit the model using Lasso Regression | 3 |
| Training and Testing Error | 5 |
| c) Fit an elastic net model on the training data | 6 |
| Report the selected tuning parameters | 6 |
| Test error | 8 |
| (d) Fit a partial least squares model | 8 |
| Fit model using plsr | 8 |
| e) Select Model | 11 |

(a) Least Square

Fit a linear model using least squares on the training data. Is there any potential disadvantage of this model?

```
house_training <- read_csv("housing_training.csv") %>%
  janitor::clean_names()
```

```
## Rows: 1440 Columns: 26
```

```
## -- Column specification -----
## Delimiter: ","
## chr (4): Overall_Qual, Kitchen_Qual, Fireplace_Qu, Exter_Qual
## dbl (22): Gr_Liv_Area, First_Flr_SF, Second_Flr_SF, Total_Bsmt_SF, Low_Qual_...
```

```
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

Perform 10-folds cross-validation on the full model

There are 25 predictors in total, we want to perform CV to find the best fitted parameters.

```
# Reproducibility
set.seed(1)
fit.lm <- train(sale_price ~ .,
               data = house_training,
               method = "lm",
               trControl = trainControl(method = "cv", number = 10))
# Print the coefficients of the final model
summary(fit.lm)
```

```
##
## Call:
## lm(formula = .outcome ~ ., data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -89864  -12424    416   12143  140205
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -4.985e+06  3.035e+06  -1.642  0.10076
## gr_liv_area     2.458e+01  1.393e+01   1.765  0.07778 .
## first_flr_sf    4.252e+01  1.409e+01   3.017  0.00260 **
## second_flr_sf   4.177e+01  1.379e+01   3.029  0.00250 **
## total_bsmt_sf   3.519e+01  2.744e+00  12.827 < 2e-16 ***
## low_qual_fin_sf          NA          NA          NA          NA
## wood_deck_sf     1.202e+01  4.861e+00   2.474  0.01350 *
## open_porch_sf    1.618e+01  1.004e+01   1.611  0.10736
## bsmt_unf_sf     -2.087e+01  1.723e+00 -12.116 < 2e-16 ***
## mas_vnr_area     1.046e+01  4.229e+00   2.473  0.01353 *
## garage_cars      4.229e+03  1.893e+03   2.234  0.02563 *
## garage_area      7.769e+00  6.497e+00   1.196  0.23195
## year_built       3.251e+02  3.130e+01  10.388 < 2e-16 ***
## tot_rms_abv_grd  -3.838e+03  6.922e+02  -5.545 3.51e-08 ***
## full_bath       -4.341e+03  1.655e+03  -2.622  0.00883 **
## overall_qualAverage -5.013e+03  1.735e+03  -2.890  0.00391 **
## overall_qualBelow_Average -1.280e+04  2.677e+03  -4.782 1.92e-06 ***
## overall_qualExcellent  7.261e+04  5.381e+03  13.494 < 2e-16 ***
## overall_qualFair     -1.115e+04  5.240e+03  -2.127  0.03356 *
## overall_qualGood      1.226e+04  1.950e+03   6.287 4.30e-10 ***
## overall_qualVery_Excellent 1.304e+05  8.803e+03  14.810 < 2e-16 ***
## overall_qualVery_Good    3.798e+04  2.741e+03  13.852 < 2e-16 ***
## kitchen_qualFair     -2.663e+04  6.325e+03  -4.210 2.71e-05 ***
## kitchen_qualGood     -1.879e+04  4.100e+03  -4.582 5.01e-06 ***
## kitchen_qualTypical   -2.677e+04  4.281e+03  -6.252 5.37e-10 ***
## fireplaces          1.138e+04  2.257e+03   5.043 5.18e-07 ***
## fireplace_quFair     -7.207e+03  6.823e+03  -1.056  0.29106
## fireplace_quGood      6.070e+02  5.833e+03   0.104  0.91713
## fireplace_quNo_Fireplace  3.394e+03  6.298e+03   0.539  0.59002
## fireplace_quPoor     -5.185e+03  7.399e+03  -0.701  0.48362
## fireplace_quTypical   -6.398e+03  5.897e+03  -1.085  0.27814
```

```
## exter_qualFair          -3.854e+04  8.383e+03  -4.598  4.66e-06 ***
## exter_qualGood         -1.994e+04  5.585e+03  -3.569  0.00037 ***
## exter_qualTypical      -2.436e+04  5.874e+03  -4.147  3.57e-05 ***
## lot_frontage           1.024e+02  1.905e+01   5.376  8.90e-08 ***
## lot_area               6.042e-01  7.864e-02   7.683  2.91e-14 ***
## longitude             -3.481e+04  2.537e+04  -1.372  0.17016
## latitude              5.874e+04  3.483e+04   1.686  0.09193 .
## misc_val              9.171e-01  1.003e+00   0.914  0.36071
## year_sold             -6.455e+02  4.606e+02  -1.401  0.16132
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 22190 on 1401 degrees of freedom
## Multiple R-squared:  0.9116, Adjusted R-squared:  0.9092
## F-statistic: 380.3 on 38 and 1401 DF,  p-value: < 2.2e-16
```

MSE, R-squared values

Check the mean squared error (MSE) and the R-squared value of the model:

```
# The Mean Squared Error is
mean((fit.lm$resample$RMSE)^2)
```

```
## [1] 531169999
```

```
# The R^2 value is
mean(fit.lm$resample$Rsquared)
```

```
## [1] 0.9029694
```

So the training error for the LS model is 531169999, with adjusted R-squared value of 0.9029694.

Potential disadvantage: 1. The model has the smallest variance among all the unbiased models, but it has larger variance compared to biased models, but it is not necessarily good enough. It may be better to have a biased estimator with smaller variance/mean squared error. 2. There are many predictors used in this model, and this can cause problem (large variance, colinearity among predictors, interpretation becomes hazardous). 3. The model is too complex, a simple model might be better. 4. This model has overfitting problem, meaning that it may perform badly on testing set.

(b) Lasso

```
x_train <- model.matrix(sale_price ~ ., house_training)[ , -1]
y_train <- house_training$sale_price
```

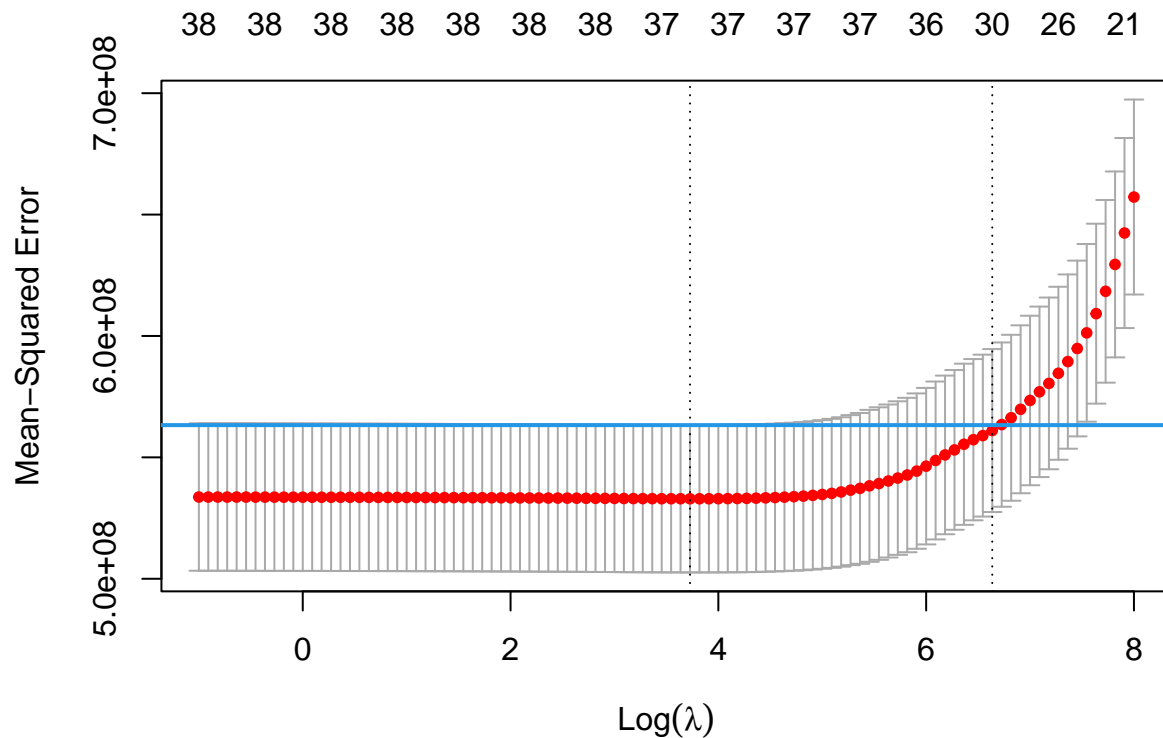
Fit the model using Lasso Regression

```

set.seed(1)
# fit the lasso regression (alpha = 1) with a sequence of lambdas
cv.lasso <- cv.glmnet(x = x_train, y = y_train,
                     standardize = TRUE,
                     alpha = 1,
                     lambda = exp(seq(8, -1, length = 100)))

plot(cv.lasso)
abline(h = (cv.lasso$cvm + cv.lasso$cvstd)[which.min(cv.lasso$cvm)], col = 4, lwd = 2)

```



When 1se rule is applied, there are total 30 predictors in the model.

The coefficient of LASSO regression model with 1SE rule applied is:

```

predict(cv.lasso, s = cv.lasso$lambda.1se, type = "coefficient")

```

```

## 40 x 1 sparse Matrix of class "dgCMatrix"
##                               s1
## (Intercept)                -2.535481e+06
## gr_liv_area                  5.723258e+01
## first_flr_sf                 1.066538e+00
## second_flr_sf                 .
## total_bsmt_sf                3.678078e+01
## low_qual_fin_sf              -2.724627e+01
## wood_deck_sf                 8.467731e+00

```

```
## open_porch_sf          8.403562e+00
## bsmt_unf_sf           -1.971379e+01
## mas_vnr_area          1.414590e+01
## garage_cars           3.078143e+03
## garage_area           1.115981e+01
## year_built            3.123467e+02
## tot_rms_abv_grd       -1.439453e+03
## full_bath              .
## overall_qualAverage   -3.158907e+03
## overall_qualBelow_Average -9.303478e+03
## overall_qualExcellent  9.066941e+04
## overall_qualFair      -6.610047e+03
## overall_qualGood       1.002368e+04
## overall_qualVery_Excellent 1.606478e+05
## overall_qualVery_Good  3.636198e+04
## kitchen_qualFair      -5.480630e+03
## kitchen_qualGood       .
## kitchen_qualTypical   -9.580041e+03
## fireplaces            6.514237e+03
## fireplace_quFair      .
## fireplace_quGood       4.630463e+03
## fireplace_quNo_Fireplace .
## fireplace_quPoor       .
## fireplace_quTypical   -3.377484e+02
## exter_qualFair        -1.481214e+04
## exter_qualGood        .
## exter_qualTypical     -5.036385e+03
## lot_frontage          7.261206e+01
## lot_area              5.648507e-01
## longitude             -1.213699e+04
## latitude              1.973031e+04
## misc_val              .
## year_sold             .
```

Training and Testing Error

Calculate the MSE of the test set

```
housing_testing <- read_csv("housing_test.csv") %>%
  janitor::clean_names()
x_test <- model.matrix(sale_price ~ ., housing_testing)[ , -1]
y_test <- housing_testing$sale_price

# Training Error
y_pred_t <- predict(cv.lasso, newx = x_train, s = "lambda.1se", type = "response")
mean(RMSE(y_pred_t, y_train)^2)

## [1] 515636300
```

The training error (MSE) is 515636300 for lasso model.

```
y_pred <- predict(cv.lasso, newx = x_test, s = "lambda.1se", type = "response")
lasso_te <- mean(RMSE(y_pred, y_test)^2)
lasso_te
```

```
## [1] 420354616
```

The test error (MSE) is 420354616 for lasso regression model when the 1SE rule is applied.

c) Fit an elastic net model on the training data

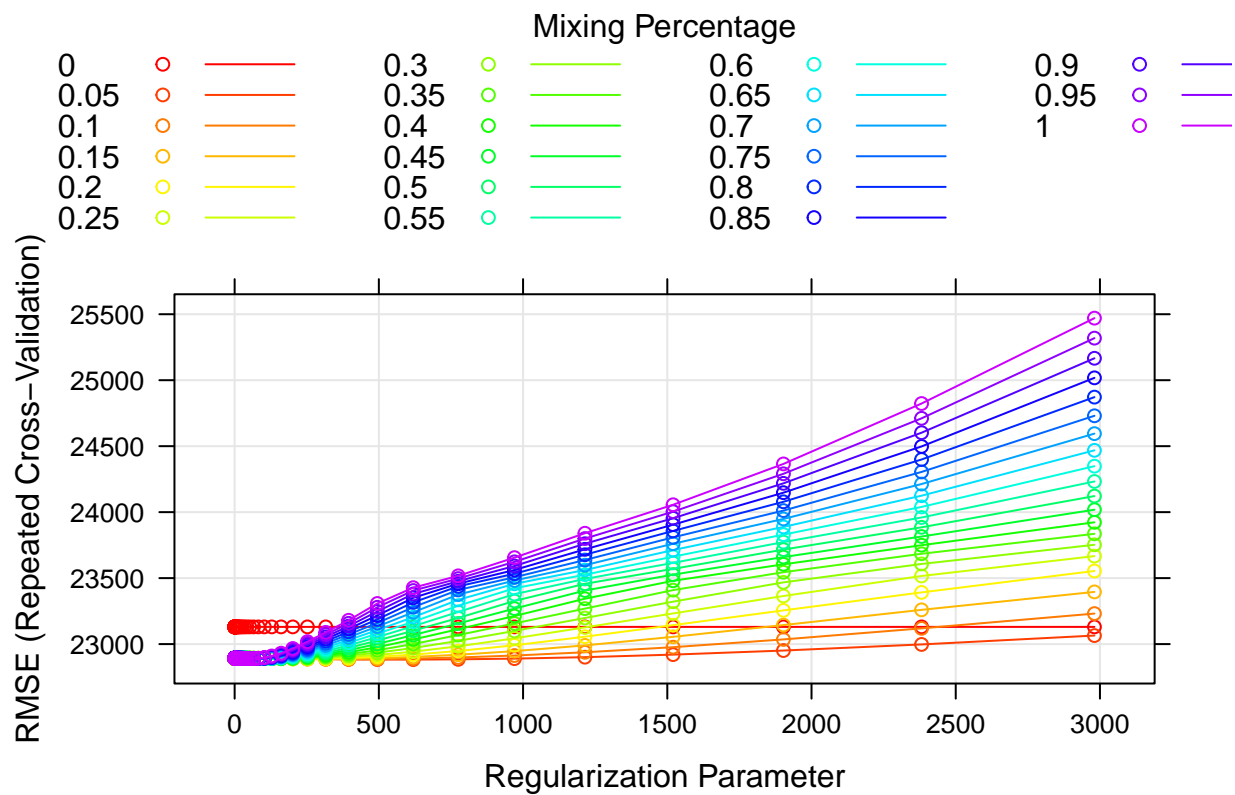
Report the selected tuning parameters

```
set.seed(2)
enet.fit <- train(x_train, y_train,
                 method = "glmnet",
                 tuneGrid = expand.grid(alpha = seq(0, 1, length = 21),
                                       lambda = exp(seq(8, -3, length = 50))),
                 trControl = trainControl(method = "repeatedcv", number = 10, repeats = 5))
enet.fit$bestTune
```

```
##      alpha      lambda
## 93  0.05 619.2886
```

```
# Set rainbow color
myCol<- rainbow(25)
myPar <- list(superpose.symbol = list(col = myCol),
             superpose.line = list(col = myCol))

plot(enet.fit, par.settings = myPar)
```



```
coef(enet.fit$finalModel, enet.fit$bestTune$lambda)
```

```
## 40 x 1 sparse Matrix of class "dgCMatrix"
##               s1
## (Intercept)  -5.119099e+06
## gr_liv_area   3.871803e+01
## first_flr_sf  2.667284e+01
## second_flr_sf 2.539870e+01
## total_bsmt_sf 3.492629e+01
## low_qual_fin_sf -1.586669e+01
## wood_deck_sf   1.234466e+01
## open_porch_sf  1.692220e+01
## bsmt_unf_sf    -2.071319e+01
## mas_vnr_area   1.172260e+01
## garage_cars    4.041408e+03
## garage_area    8.948237e+00
## year_built     3.188394e+02
## tot_rms_abv_grd -3.419303e+03
## full_bath      -3.659651e+03
## overall_qualAverage -5.119003e+03
## overall_qualBelow_Average -1.269888e+04
## overall_qualExcellent 7.592629e+04
## overall_qualFair -1.148475e+04
## overall_qualGood 1.195840e+04
## overall_qualVery_Excellent 1.365876e+05
```

```
## overall_qualVery_Good      3.762260e+04
## kitchen_qualFair          -2.354051e+04
## kitchen_qualGood          -1.597423e+04
## kitchen_qualTypical       -2.402780e+04
## fireplaces                 1.079788e+04
## fireplace_quFair          -7.863526e+03
## fireplace_quGood           1.474300e+02
## fireplace_quNo_Fireplace   1.757694e+03
## fireplace_quPoor          -5.809405e+03
## fireplace_quTypical       -6.964071e+03
## exter_qualFair            -3.276755e+04
## exter_qualGood            -1.436584e+04
## exter_qualTypical         -1.897822e+04
## lot_frontage               1.000357e+02
## lot_area                   6.030183e-01
## longitude                  -3.517964e+04
## latitude                   5.769405e+04
## misc_val                   8.652364e-01
## year_sold                  -5.712351e+02
```

The tuning parameter λ is 619.2886. The parameter α is 0.05.

Test error

```
enet.pred <- predict(enet.fit, newdata = x_test)
# test error
enet_te <- mean(RMSE(enet.pred, y_test)^2)
enet_te
```

```
## [1] 438209306
```

The test error (MSE) is 438209306 for elastic net model.

(d) Fit a partial least squares model

Fit model using plsr

```
set.seed(2)
pls.mod <- plsr(sale_price~.,
                data = house_training,
                scale = TRUE,
                validation = "CV")

summary(pls.mod)
```

```
## Data:      X dimension: 1440 39
## Y dimension: 1440 1
## Fit method: kernelpls
```

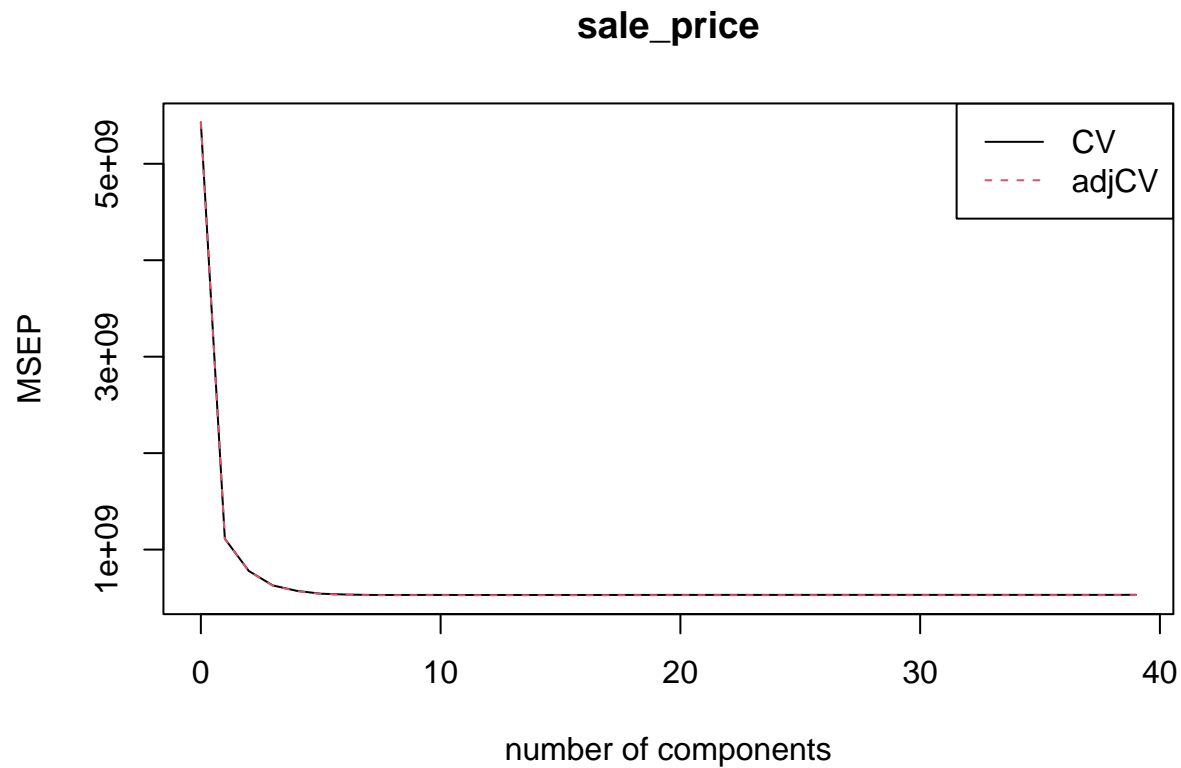


```

## Number of components considered: 39
##
## VALIDATION: RMSEP
## Cross-validated using 10 random segments.
##      (Intercept)  1 comps  2 comps  3 comps  4 comps  5 comps  6 comps
## CV              73685   33357   27895   25075   23920   23294   23114
## adjCV           73685   33354   27863   25006   23858   23233   23061
##      7 comps  8 comps  9 comps 10 comps 11 comps 12 comps 13 comps
## CV      23020   23007   23021   23022   23011   23011   23013
## adjCV    22968   22954   22965   22964   22953   22952   22953
##      14 comps 15 comps 16 comps 17 comps 18 comps 19 comps 20 comps
## CV      23010   23016   23014   23014   23019   23022   23024
## adjCV    22951   22956   22954   22954   22959   22961   22964
##      21 comps 22 comps 23 comps 24 comps 25 comps 26 comps 27 comps
## CV      23027   23031   23029   23030   23030   23032   23033
## adjCV    22967   22970   22968   22969   22969   22971   22972
##      28 comps 29 comps 30 comps 31 comps 32 comps 33 comps 34 comps
## CV      23034   23034   23034   23034   23034   23034   23034
## adjCV    22973   22973   22972   22973   22973   22973   22973
##      35 comps 36 comps 37 comps 38 comps 39 comps
## CV      23034   23034   23034   23034   23052
## adjCV    22973   22973   22973   22973   22990
##
## TRAINING: % variance explained
##      1 comps  2 comps  3 comps  4 comps  5 comps  6 comps  7 comps
## X              20.02   25.93   29.67   33.59   37.01   40.03   42.49
## sale_price     79.73   86.35   89.36   90.37   90.87   90.99   91.06
##      8 comps  9 comps 10 comps 11 comps 12 comps 13 comps 14 comps
## X              45.53   47.97   50.15   52.01   53.69   55.35   56.86
## sale_price     91.08   91.10   91.13   91.15   91.15   91.16   91.16
##      15 comps 16 comps 17 comps 18 comps 19 comps 20 comps
## X              58.64   60.01   62.18   63.87   65.26   67.10
## sale_price     91.16   91.16   91.16   91.16   91.16   91.16
##      21 comps 22 comps 23 comps 24 comps 25 comps 26 comps
## X              68.44   70.12   71.72   73.35   75.20   77.27
## sale_price     91.16   91.16   91.16   91.16   91.16   91.16
##      27 comps 28 comps 29 comps 30 comps 31 comps 32 comps
## X              78.97   80.10   81.83   83.55   84.39   86.34
## sale_price     91.16   91.16   91.16   91.16   91.16   91.16
##      33 comps 34 comps 35 comps 36 comps 37 comps 38 comps
## X              88.63   90.79   92.79   95.45   97.49   100.00
## sale_price     91.16   91.16   91.16   91.16   91.16   91.16
##      39 comps
## X              100.67
## sale_price     91.16

```

```
validationplot(pls.mod, val.type="MSEP", legendpos = "topright")
```



Training Error and Testing Error

```
# training error
cv.mse <- RMSEP(pls.mod)
mean(min(cv.mse$val[1,,])^2)
```

[1] 529300397

```
ncomp.cv <- which.min(cv.mse$val[1,,])-1 # extract the response and delete the 0th component
# num of components
ncomp.cv
```

```
## 8 comps
##      8
```

```
# Prediction
pls_pred <- predict(pls.mod, newdata = x_test,
                    ncomp = ncomp.cv)

# test MSE
pls_te <- mean(RMSE(y_test, pls_pred)^2)
pls_te
```

[1] 440217938

The model with 8 components yields the minimum training error (MSE) 529300397. The testing error (MSE) is 440217938.

e) Select Model

```
#compute the testing error of LS regression first
y_pred <- predict(fit.lm, newdata = housing_testing)
ls_te <- mean(RMSE(y_pred, y_test)^2)
ls_te
```

```
## [1] 447287652
```

The test error of LS model is 447287652.

For response predicting, we want to choose the model with the smallest test error. The test error (MSE) is 420354616 for lasso regression model, 438209306 for elastic net model, and 440217938 for the PLS model. Therefore, we should choose the lasso regression model for predicting the response.