

## Exam 1 Advance R

```
setwd("C:/Users/akash/Desktop/Spring 2017/Advance R/Exam1")
getwd()
```

```
## [1] "C:/Users/akash/Desktop/Spring 2017/Advance R/Exam1"
```

1. Load the dataset into R.

```
mydata <- read.csv("exam1.csv", header=TRUE, sep=",")
```

2. Center and scale numerical predictors.

```
sc <- scale(mydata[sapply(mydata, is.numeric)], center = TRUE, scale = TRUE)
sc <- cbind.data.frame(as.data.frame(sc), as.data.frame(mydata$X15))
head(sc)
```

```
##           X0           X1           X2           X3           X4           X5
## 1  1.38162065 -0.5852751  1.3577944 -0.5975150  0.5703694 -0.06469633
## 2 -1.56241660 -0.3580871  0.6811227 -0.4719006  0.7505207  0.14539129
## 3 -0.19884214  0.2703088 -1.4210064  0.8028125  1.1659734 -1.73838751
## 4 -1.16158464 -0.6221013  1.0742198 -1.3521837 -0.7687214 -0.54758265
## 5 -0.05927866 -1.0997133 -0.4569793 -1.6009201 -0.1422798  0.91642468
## 6  0.11180996 -0.4569426 -2.9033014  0.3662817 -0.3441083  1.35518753
##           X6           X7           X8           X9           X11           X12
## 1 -0.2192142  1.332037606  2.0881376 -1.768388062 -0.7795136 -0.38284320
## 2  2.1730191  1.521699533 -1.5106584 -1.155575951 -1.2099636 -1.51853753
## 3  0.9294798 -0.002486363 -0.4190991 -0.602812834 -0.3757958 -0.06513682
## 4  0.2611978 -0.457903838 -0.5847360  0.735045193  0.2404122  0.44265588
## 5  0.6185339 -1.070999417 -0.4020839 -0.004904503 -0.7362294  0.04865450
## 6  0.2487695  0.620745584 -0.9163592 -0.259813570 -1.3366131  0.26361913
##           X13           X14           y mydata$X15
## 1  0.41282067 -1.39453397 -0.86391656      True
## 2  0.30790277  1.13590230 -0.70634478     False
## 3 -0.29228272  0.83857008 -0.34671627      True
## 4  0.33684205 -0.23564971 -0.30349688      True
## 5 -0.79376489  1.04065697 -0.03099074      True
## 6 -0.07521933  0.01437285 -0.91015485      True
```

```
attr(,"scaled:center")
```

```
## NULL
```

```
attr(,"scaled:scale")
```

```
## NULL
```

3. Create dummy variables for any categorical predictors.

```
dX15 <- as.numeric(mydata$X15 == "True")
sc <- cbind.data.frame(as.data.frame(sc), as.data.frame(dX15))
```

```
colnames(sc)[colnames(sc)=="mydata$X15"] <- "X15"
head(sc)
```

```
##           X0           X1           X2           X3           X4           X5
## 1  1.38162065 -0.5852751  1.3577944 -0.5975150  0.5703694 -0.06469633
## 2 -1.56241660 -0.3580871  0.6811227 -0.4719006  0.7505207  0.14539129
## 3 -0.19884214  0.2703088 -1.4210064  0.8028125  1.1659734 -1.73838751
## 4 -1.16158464 -0.6221013  1.0742198 -1.3521837 -0.7687214 -0.54758265
## 5 -0.05927866 -1.0997133 -0.4569793 -1.6009201 -0.1422798  0.91642468
## 6  0.11180996 -0.4569426 -2.9033014  0.3662817 -0.3441083  1.35518753
##           X6           X7           X8           X9           X11           X12
## 1 -0.2192142  1.332037606  2.0881376 -1.768388062 -0.7795136 -0.38284320
## 2  2.1730191  1.521699533 -1.5106584 -1.155575951 -1.2099636 -1.51853753
## 3  0.9294798 -0.002486363 -0.4190991 -0.602812834 -0.3757958 -0.06513682
## 4  0.2611978 -0.457903838 -0.5847360  0.735045193  0.2404122  0.44265588
## 5  0.6185339 -1.070999417 -0.4020839 -0.004904503 -0.7362294  0.04865450
## 6  0.2487695  0.620745584 -0.9163592 -0.259813570 -1.3366131  0.26361913
##           X13           X14           y           X15  dx15
## 1  0.41282067 -1.39453397 -0.86391656  True      1
## 2  0.30790277  1.13590230 -0.70634478 False     0
## 3 -0.29228272  0.83857008 -0.34671627  True      1
## 4  0.33684205 -0.23564971 -0.30349688  True      1
## 5 -0.79376489  1.04065697 -0.03099074  True      1
## 6 -0.07521933  0.01437285 -0.91015485  True      1
```

4. Split the data into a training and test set. Set aside the test set until the end.

```
sSize = nrow(sc)*0.8
Splitteddata <- sample(nrow(sc), sSize, replace = TRUE)
trainData<- sc[Splitteddata,]
head(trainData)
```

```
##           X0           X1           X2           X3           X4           X5
## 44 -1.6568179 -2.4419388  0.9879263 -0.15980442  0.09401606  1.12260380
## 92 -0.9013076 -1.0646442 -0.3879080  1.76440520 -1.32730324  0.09847768
## 100 -1.1287533  1.4557771  1.6544679 -0.37677395 -0.34701597 -0.19153968
## 75 -0.6008111 -2.7256421 -1.7091325  0.03735086  0.65207879  0.43114010
## 8   0.3297880  2.4150009  0.4053096  0.63536367 -0.11381102  0.13477983
## 62  1.0510132  0.3329312 -0.7070745 -0.54772484 -2.92662622 -0.50400401
##           X6           X7           X8           X9           X11           X12
## 44 -1.9010534 -0.71697953 -1.5053475 -0.3620340  0.4728863  0.922121544
## 92 -0.4635289  0.73056830  2.0507000 -1.1164068  1.6381708 -1.711545323
## 100 -0.7593635  0.10733795 -0.6386472  0.4819096 -1.9586547 -0.719446203
## 75 -1.5828404  1.05341242 -0.8249170 -0.3019082  0.3535312  2.178729008
## 8   0.5559363  0.48713237 -1.0046308 -0.4833098  0.7773271  0.819050700
## 62  1.9992188  0.03332573  0.2040490  0.5804240 -0.2238018 -0.001399971
##           X13           X14           y           X15  dx15
## 44 -0.24271776  1.3167526  0.18621716  True      1
## 92  0.05723332 -0.2662174 -1.38107230  True      1
## 100 -0.70507342  0.1850330  0.15133607 False     0
## 75  0.96914965  1.3607625 -0.09846648 False     0
```

```
## 8      0.28940416  0.6271701  0.77350224 False    0
## 62     -0.66526530 -1.0535356 -0.42447685 False    0
```

```
testData<- sc[-Splitteddata,]
head(testData)
```

```
##           X0           X1           X2           X3           X4           X5
## 1  1.381620645 -0.5852751  1.3577944 -0.5975150  0.57036938 -0.06469633
## 11 -0.173504877  0.3703515  1.1412113 -0.1894249 -1.06465209 -0.07561962
## 12 -0.009359255  0.4682381  0.6772946  0.3938466  0.39750909 -1.91753819
## 13 -0.220018949 -2.3301803  1.1195224 -0.7477318  0.09952701  0.12311452
## 15  1.169910457 -0.3786583  0.2429107 -0.8905387 -0.50907246  1.24675172
## 18  1.071014710  1.0010285 -0.9595237 -0.7671195 -0.73899843  0.34011330
##           X6           X7           X8           X9           X11           X12
## 1  -0.2192142  1.33203761  2.0881376 -1.7683881 -0.7795136 -0.38284320
## 11 -1.0475209  0.44778365  0.8564377  0.2262502 -2.0090915  0.30950336
## 12 -0.9573937  0.06440245 -1.4850723  0.2233342 -0.3788495  0.78787297
## 13 -0.3421736  0.49150600  0.9353150  0.4746702  0.5994617 -0.56730478
## 15 -0.5105489 -0.00209826  0.4030609  1.2384760  1.1930902 -1.50155333
## 18 -1.7880368 -1.32549756 -1.7969893  0.3072174 -0.6407095  0.09440536
##           X13           X14           y           X15  dX15
## 1  0.412820671 -1.3945340 -0.8639166  True    1
## 11 -0.004587436  0.4909828  0.5284274 False    0
## 12  1.562377786  1.3779212  0.9271684  True    1
## 13 -0.735211359  1.0080568  0.7821590  True    1
## 15 -0.167721709  0.7591709  1.1585631  True    1
## 18  1.277279599 -1.0793178 -0.9729553 False    0
```

5. Split the training data using 4 fold cross validation.

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.3.2
```

```
## Loading required package: lattice
```

```
## Loading required package: ggplot2
```

```
trainDataFolds <- createFolds(trainData$X15, k = 4, list = TRUE)
```

6. Fit ridge regression models for a range of " $\lambda^2$ " values. Be sure to include large enough values of " $\lambda^2$ " that you see a decrease in performance.

7. For each value of " $\lambda^2$ ", you will have 4 models (1 for each fold). Evaluate the RMSE of all models on the fold not used to train. Use a loop for this.

```
library(AppliedPredictiveModeling)
```

```
## Warning: package 'AppliedPredictiveModeling' was built under R version
## 3.3.2
```

```
library(penalized)
```

```

## Warning: package 'penalized' was built under R version 3.3.2
## Loading required package: survival
##
## Attaching package: 'survival'
##
## The following object is masked from 'package:caret':
##
##      cluster
## Welcome to penalized. For extended examples, see vignette("penalized").

library(glmnet)

## Warning: package 'glmnet' was built under R version 3.3.2
## Loading required package: Matrix
## Loading required package: foreach
## Loaded glmnet 2.0-5

library(elasticnet)

## Warning: package 'elasticnet' was built under R version 3.3.2
## Loading required package: lars
## Warning: package 'lars' was built under R version 3.3.2
## Loaded lars 1.2

library(MASS)

## Warning: package 'MASS' was built under R version 3.3.2

rmse <- function(a,b){
  x= (a-b)^2
  y= sqrt(mean(x))
  return(y)
}

perf.df <- data.frame("q"=numeric(0))
average.root <- data.frame("t"= numeric(0))
sequence <- seq(0, 1,0.01 )
for (i in 1:length(sequence)) {
  for (j in length(trainDataFolds)) {
    crossvalidationTest <- trainData[trainDataFolds[[j]],]
    crossvalidationTrain <- trainData[-trainDataFolds[[j]],]
    crossvalidationTestCopy <- crossvalidationTest
    crossvalidationTest$X15 <- NULL
    crossvalidationTest$y <- NULL
    ridge.reg <- penalized(y ~

```

[illegible]

[illegible]

## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12  
## 12

perf.df

##	X0.710231686603127
## 1	0.7102317
## 2	0.7102317
## 3	0.7102318
## 4	0.7102318
## 5	0.7102319
## 6	0.7102320
## 7	0.7102321
## 8	0.7102322
## 9	0.7102324
## 10	0.7102325
## 11	0.7102327
## 12	0.7102328
## 13	0.7102330
## 14	0.7102332
## 15	0.7102334
## 16	0.7102336
## 17	0.7102339
## 18	0.7102341
## 19	0.7102343
## 20	0.7102346
## 21	0.7102349
## 22	0.7102352
## 23	0.7102355
## 24	0.7102358
## 25	0.7102361
## 26	0.7102365
## 27	0.7102368
## 28	0.7102372
## 29	0.7102376
## 30	0.7102380
## 31	0.7102384
## 32	0.7102388
## 33	0.7102392

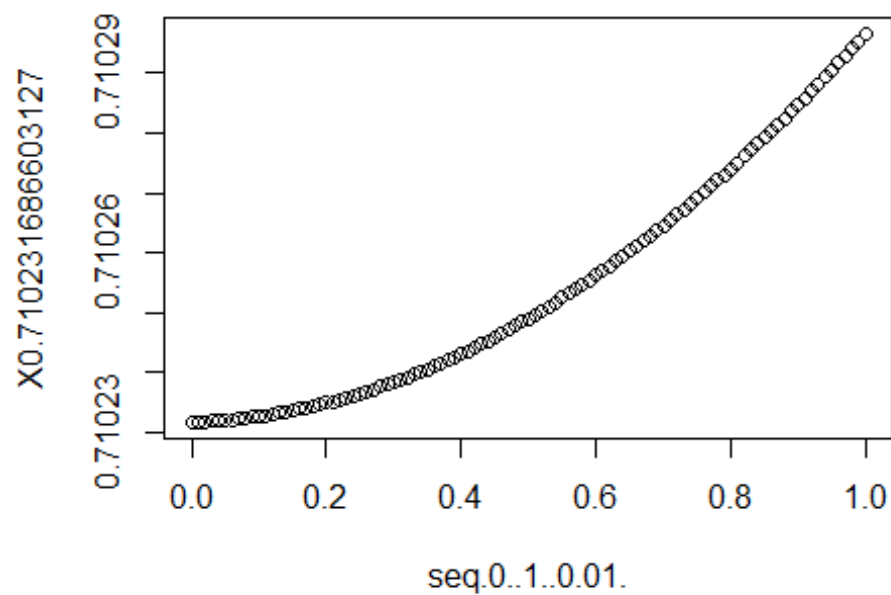
## 34	0.7102397
## 35	0.7102401
## 36	0.7102406
## 37	0.7102411
## 38	0.7102416
## 39	0.7102421
## 40	0.7102426
## 41	0.7102431
## 42	0.7102436
## 43	0.7102442
## 44	0.7102448
## 45	0.7102453
## 46	0.7102459
## 47	0.7102465
## 48	0.7102471
## 49	0.7102478
## 50	0.7102484
## 51	0.7102490
## 52	0.7102497
## 53	0.7102504
## 54	0.7102511
## 55	0.7102518
## 56	0.7102525
## 57	0.7102532
## 58	0.7102539
## 59	0.7102547
## 60	0.7102554
## 61	0.7102562
## 62	0.7102570
## 63	0.7102578
## 64	0.7102586
## 65	0.7102594
## 66	0.7102602
## 67	0.7102611
## 68	0.7102619
## 69	0.7102628
## 70	0.7102637
## 71	0.7102645
## 72	0.7102654
## 73	0.7102664
## 74	0.7102673
## 75	0.7102682
## 76	0.7102692
## 77	0.7102701
## 78	0.7102711
## 79	0.7102721
## 80	0.7102731
## 81	0.7102741
## 82	0.7102751
## 83	0.7102761



```
## 84      0.7102771
## 85      0.7102782
## 86      0.7102793
## 87      0.7102803
## 88      0.7102814
## 89      0.7102825
## 90      0.7102836
## 91      0.7102848
## 92      0.7102859
## 93      0.7102870
## 94      0.7102882
## 95      0.7102893
## 96      0.7102905
## 97      0.7102917
## 98      0.7102929
## 99      0.7102941
## 100     0.7102954
## 101     0.7102966
```

8. Make a plot with " $\lambda^2$ " on the x-axis and the mean RMSE (average over the 4 folds) on the y-axis.

```
rmse.plot <- data.frame(seq(0, 1, 0.01), perf.df)
plot(rmse.plot)
```



9. Using this plot, select "lambda"2 for your model. Explain your reasoning. From the plot it can be observed that the value of root mean square error first increases as lambda increases and then decreases. It is minimum near to 0.4. Since we want RMSE to be minimized to get a good model, So, the value of lambda2 can be taken as 0.4.

10. Fit a model on the complete training data using your selected value for "lambda"2.

```
Final.regression <- penalized(y ~  
X0+X1+X2+X3+X4+X5+X6+X7+X8+X9+X11+X12+X13+X14+dX15, data = trainData ,  
lambda2 = 0.4, standardize = TRUE)
```

```
## 12
```

```
Final.regression
```

```
## Penalized linear regression object
```

```
## 16 regression coefficients
```

```
##
```

```
## Loglikelihood = 156.9725
```

```
## L2 penalty = 0.2257771 at lambda2 = 0.4
```

11. Evaluate the R2 and RMSE of your model on the test set.

```
testDataCopy <- testData
```

```
testData$X15 <- NULL
```

```
testData$y <- NULL
```

```
Final.predict <- predict(Final.regression, testData)
```

```
Final.rmse <- rmse(Final.predict, testDataCopy$y)
```

```
Final.rmse
```

```
## [1] 0.7368192
```

```
mean.testdata <- lapply(testDataCopy$y, mean)
```

```
SumofSquares.Total = sum((testDataCopy$y - mean(testDataCopy$y))^2)
```

```
sumofSquares.Errors = sum((testDataCopy$y - Final.predict[ , "mu"])^2)
```

```
Rsquared = 1-(sumofSquares.Errors/SumofSquares.Total)
```

```
Rsquared
```

```
## [1] 0.9982672
```