

Analytics Modeling_HW8

Fall 2024

```
# Load data
crime_df <- read.table("uscrime.txt", stringsAsFactors = FALSE, header = TRUE)
tail(crime_df, 5)

##      H  So  Ed  Poi  Po2  LF  H.F  Pop  NM  U1  U2  Wealt  Ineq  Prob
## 43 16.2  1  9.8  7.5  7.0 0.532  99.46  40 20.8 0.073 2.7  4960 22.4 0.054902
## 44 13.6  0 12.1  9.5  9.6 0.574 101.2  29  3.6 0.111 2.7  6220 16.2 0.028100
## 45 13.9  1  8.8  4.6  4.1 0.480  96.8  19  4.9 0.135 5.3  4570 24.9 0.056202
## 46 12.0  0 10.4 10.6  9.7 0.599  98.9  40  2.4 0.078 2.5  5930 17.1 0.046598
## 47 13.0  0 12.1  9.0  9.1 0.623 104.9  3  2.2 0.113 4.0  5880 16.0 0.052802

##      1 time Crime
## 43 10.9909  823
## 44 10.0001 1030
## 45 32.5996  455
## 46 16.6999  508
## 47 16.0997  849
```

Question 11.1

Using the crime data set, build a regression model using Stepwise Regression

I started with a model that used all predictors and then used stepwise regression to reduce the number of variables. At each step, the stepwise regression removed the predictor with the lowest AIC until finally the dataset was reduced from 15 predictors to 8.

The initial model with all predictors resulted in an adjusted R-squared value of 0.7078, indicating that 70.78% of the model's variability is explained by the predictors.

The stepwise regression resulted in a model containing 8 predictors:

- H
- Ed
- Poi
- H.F
- U1
- U2
- Ineq
- Prob

And resulted in an adjusted R-squared value of 0.7444, indicating that 74.44% of the model's variability is explained by the predictors.

Stepwise regression resulted in a model that improved in two ways:

- improved evaluation metrics
- increased model simplicity

```
# set seed for reproducibility
set.seed(1)

# start with a model that has all predictors
initial_model <- lm(Crime~., data = crime_df)

summary(initial_model)

##
## Call:
## lm(formula = Crime ~ ., data = crime_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -395.74  -98.09   -6.69   112.99   512.67
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -5.984e+03  1.626e+03  -3.675 0.000993 ***
## H             0.783e+01  4.171e+01   0.186 0.045443 *
## So            -3.803e+00  1.488e+02  -0.026 0.979765
## Ed             1.883e+02  6.209e+01   3.033 0.004861 **
## Poi            1.928e+02  1.061e+02   1.817 0.078892 .
## Po2           -1.094e+02  1.175e+02  -0.931 0.358830
## LF            -0.630e+02  1.470e+03  -0.452 0.654554
## H.F           1.741e+01  2.035e+01   0.855 0.390995
## Pop           -7.330e-01  1.290e+00  -0.568 0.573845
## NM            4.204e+00  6.481e+00   0.649 0.521279
## U1            -5.827e+03  4.210e+03  -1.384 0.176238
## U2            1.678e+02  8.234e+01   2.036 0.050161 .
## Wealt         0.617e+02  1.037e+01   0.928 0.360754
## Ineq          7.067e+01  2.272e+01   3.111 0.003983 **
## Prob         -4.855e+03  2.272e+03  -2.137 0.040627 *
## Time         -3.479e+00  7.165e+00  -0.486 0.630708
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 209.1 on 31 degrees of freedom
## Multiple R-squared:  0.8031, Adjusted R-squared:  0.7078
## F-statistic: 8.429 on 15 and 31 DF, p-value: 3.539e-07
```

```
# perform both direction stepwise regression
stepwise_model <- step(initial_model
                        , scope = list(lower = formula(lm(Crime~1, data = crime_df))
                                       , upper = formula(lm(Crime~., data = crime_df))
                                       )
                        , direction = "both"
                        )
```

```
## Start: AIC=514.65
## Crime ~ H + So + Ed + Poi + Po2 + LF + H.F + Pop + NM + U1 +
##           U2 + Wealt + Ineq + Prob + Time
##
##              Df Sum of Sq  RSS   AIC
## - So           1       29 1354974 512.65
## - LF           1      8917 1363862 512.96
## - Time         1     10304 1365250 513.00
## - Pop          1     14122 1369068 513.14
## - NM           1     10395 1373401 513.28
## - H.F          1     31967 1386913 513.74
## - Wealt        1     37613 1392558 513.94
## - Po2          1     37919 1392865 513.95
## <none>                 1354946 514.65
## - U1           1     83722 1438668 515.47
## - Poi          1     144306 1449922 517.41
## - U2           1     181536 1516482 516.81
## - M            1     193770 1548716 518.93
## - Prob         1     199538 1554484 519.11
## - Ed           1     402117 1757063 524.86
## - Ineq         1     423031 1777977 525.42

## Step: AIC=512.65
## Crime ~ H + Ed + Poi + Po2 + LF + H.F + Pop + NM + U1 + U2 +
##           Wealt + Ineq + Prob + Time
##
##              Df Sum of Sq  RSS   AIC
## - Time         1     10341 13565315 511.03
## - LF           1     10878 1365852 511.03
## - Pop          1     14127 1369101 511.14
## - NM           1     21626 1376600 511.39
## - H.F          1     32449 1387423 511.76
## - Po2          1     37954 1392939 511.95
## - Wealt        1     39223 1394197 511.99
## <none>                 1354974 512.65
## - U1           1     96420 1451395 513.88
## - So           1     29 1354946 514.65
## - Poi          1     144306 1449922 517.41
## - U2           1     189859 1544834 516.81
## - M            1     195084 1559059 516.97
## - Prob         1     204463 1559437 517.26
## - Ed           1     403140 1758114 522.89
## - Ineq         1     488834 1843808 525.13

## Step: AIC=511.01
## Crime ~ H + Ed + Poi + Po2 + LF + H.F + Pop + NM + U1 + U2 +
##           Wealt + Ineq + Prob
##
##              Df Sum of Sq  RSS   AIC
## - LF           1     10533 1375848 509.37
## - NM           1     15482 1380797 509.54
## - Pop          1     21846 1387161 509.75
## - Po2          1     28932 1394247 509.99
## - Wealt        1     36070 1401385 510.23
## - H.F          1     41784 1407990 510.42
## <none>                 1365315 511.01
## - U1           1     91420 1456735 512.05
## - Time         1     10341 1354974 512.65
## - So           1     65 1365250 513.00
## - Poi          1     13437 1499452 513.41
## - U2           1     184143 1549456 514.95
## - M            1     186110 1551425 515.01
## - Prob         1     237493 1602808 516.54
## - Ed           1     409448 1774763 521.33
## - Ineq         1     582909 1868224 523.75

## Step: AIC=509.37
## Crime ~ H + Ed + Poi + Po2 + LF + H.F + Pop + NM + U1 + U2 + Wealt +
##           Ineq + Prob
##
##              Df Sum of Sq  RSS   AIC
## - NM           1     11675 1397523 507.77
## - Po2          1     21418 1397266 508.09
## - Pop          1     27803 1403651 508.31
## - H.F          1     31252 1407100 508.42
## - Wealt        1     35035 1410883 508.55
## <none>                 1375848 509.37
## - U1           1     80954 1456802 510.06
## - LF           1     10533 1365315 511.01
## - Time         1     9996 1365852 511.03
## - So           1     3046 1372802 511.26
## - Poi          1     123896 1499744 511.42
## - U2           1     190746 1566954 513.47
## - M            1     217716 1593564 514.27
## - Prob         1     226971 1602819 514.54
## - Ed           1     413254 1789103 519.71
## - Ineq         1     500944 1876792 521.96

## Step: AIC=507.77
## Crime ~ H + Ed + Poi + Po2 + LF + H.F + Pop + U1 + U2 + Wealt + Ineq +
##           Prob
##
##              Df Sum of Sq  RSS   AIC
## - Po2          1     16796 1404229 506.39
## - Pop          1     25793 1413315 506.63
## - H.F          1     26785 1414308 506.66
## - Wealt        1     31551 1419073 506.82
## <none>                 1387523 507.77
## - U1           1     83881 1471404 508.52
## - NM           1     11675 1375848 509.37
## - So           1     7207 1380116 509.52
## - LF           1     6726 1380797 509.54
## - Time         1     4534 1382989 509.61
## - Poi          1     118348 1505871 509.61
## - U2           1     201453 1588976 512.14
## - Prob         1     216760 1604262 512.59
## - M            1     309214 1696737 515.22
## - Ed           1     402754 1790276 517.74
## - Ineq         1     589736 1977259 522.41

## Step: AIC=506.33
## Crime ~ H + Ed + Poi + H.F + Pop + U1 + U2 + Wealt + Ineq +
##           Prob
##
##              Df Sum of Sq  RSS   AIC
## - Pop          1     22345 1426575 505.07
## - Wealt        1     32142 1436371 505.39
## - H.F          1     36808 1441037 505.54
## <none>                 1404247 506.33
## - U1           1     86373 1490602 507.13
## - Po2          1     16706 1387523 507.77
## - NM           1     6963 1397266 508.09
## - So           1     3807 1400812 508.20
## - LF           1     1986 1402243 508.26
## - Time         1     575 1403654 508.31
## - U2           1     205814 1610043 510.76
## - Prob         1     218607 1622836 511.13
## - M            1     307061 1711230 513.62
## - Ed           1     389502 1793731 515.83
## - Ineq         1     608627 2012856 521.25
## - Poi          1     1050202 2454432 530.57

## Step: AIC=505.07
## Crime ~ H + Ed + Poi + H.F + U1 + U2 + Wealt + Ineq + Prob
##
##              Df Sum of Sq  RSS   AIC
## <none>                 1426575 505.07
## - H.F          1     84499 1511805 505.77
## - U1           1     89463 1526037 506.24
## - Pop          1     22345 1404229 506.33
## - Po2          1     13259 1413315 506.63
## - NM           1     5927 1420648 506.87
## - So           1     5724 1420851 506.88
## - LF           1     5176 1422106 506.90
## - Time         1     3913 1422661 506.94
## - Prob         1     198571 1625145 509.20
## - U2           1     208880 1635455 509.49
## - M            1     320926 1747501 512.61
## - Ed           1     386775 1815348 514.35
## - Ineq         1     594779 2021354 519.45
## - Poi          1     1127277 2553552 530.44

## Step: AIC=503.93
## Crime ~ H + Ed + Poi + H.F + U1 + U2 + Ineq + Prob
##
##              Df Sum of Sq  RSS   AIC
## <none>                 1453068 503.93
## - Wealt        1     26493 1426575 505.07
## - H.F          1     103159 1556227 505.16
## - Pop          1     103159 1556227 505.16
## - Po2          1     14148 1438919 505.47
## - So           1     9329 1443739 505.63
## - LF           1     4374 1448694 505.81
## - NM           1     3799 1449269 505.81
## - Time         1     2293 1450775 505.86
## - U2           1     127044 1540712 505.87
## - Prob         1     247978 1701846 509.34
## - U1           1     255443 1708511 509.55
## - M            1     296790 1749858 510.67
## - Ed           1     445788 1898855 514.51
## - Ineq         1     738244 2191312 521.24
## - Poi          1     1672038 3125105 537.93
```

```
summary(stepwise_model)
```

```
##
## Call:
## lm(formula = Crime ~ H + Ed + Poi + H.F + U1 + U2 + Ineq + Prob,
##     data = crime_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -444.70 -111.07   3.03  122.15  483.30
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -6426.10     1194.61  -5.379 4.04e-06 ***
## H              93.32         33.50   2.786 0.00828 **
## Ed           180.12         52.75   3.414 0.00153 **
## Poi          102.65         15.52   6.613 8.26e-08 ***
## H.F           22.34         13.60   1.642 0.10874
## U1          -6086.36     3339.27  -1.823 0.07622 .
## U2           187.35         72.48   2.585 0.01371 *
## Ineq         61.33         13.96   4.394 6.63e-05 ***
## Prob        -3796.03         1490.65  -2.547 0.01505 *
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 195.5 on 38 degrees of freedom
## Multiple R-squared:  0.7888, Adjusted R-squared:  0.7444
## F-statistic: 17.74 on 8 and 38 DF, p-value: 1.159e-10
```

Using the crime data set, build a regression model using Lasso Regression

I created a lasso model using `cv.glmnet`, which automatically scaled the data for me. The tau / lambda threshold that worked best for the model was 8.33952725, which best minimized the MSE, and resulted in a regression with 11 non-zero variables. In other words, the variable selection process removed 4 predictors. The most that best minimized the MSE resulted in an R² of 0.7743174.

```
# set seed for reproducibility
set.seed(1)

# predictors and response
X = as.matrix(crime_df[,16])
y = as.matrix(crime_df[,16])

# do k-fold cross validation for lasso model
lasso_model <- cv.glmnet(x = X
                        , y = y
                        , alpha = 1 # lasso regression alpha = 1
                        , nfolds = 8 # number of folds
                        , nlambda = 20 # tau thresholds randomly generated
                        , type.measure = "mse" # squared error for gaussian models
                        , family = "gaussian"
                        , standardize = TRUE # use automatically scaled data
                        )

# plot MSE of lasso model
plot(lasso_model)
```



```
# use the tau / lambda that corresponds to the lowest MSE
lasso_model$lambda.min
```

```
## [1] 8.339527
```

```
# get a list of tau/lambda, cross-validation error, and number of non-zero coefficients for each lambda.
cbind(lasso_model$lambda, lasso_model$cvm, lasso_model$zero)
```

```
##              [,1]      [,2] [,3]
## #0 263.09593604 151408.06  0
## #1 324.05365872 126536.67  2
## #2 162.02682936 121743.71  1
## #3 99.78393301 103954.55  1
## #4 61.45175663 93168.04  4
## #5 37.84494539 76653.63  5
## #6 14.35341873 68765.92  10
## #7 8.83952725 67508.55  11
## #8 5.44380704 67649.74  12
## #9 3.35255883 68280.67  12
## #10 2.06466736 70005.68  13
## #11 1.27352170 71798.41  14
## #12 0.70306456 73344.38  15
## #13 0.48224879 73634.73  15
## #14 0.29699205 73757.71  15
## #15 0.18290282 73917.56  15
## #16 0.11263988 74034.73  14
## #17 0.06936987 74151.66  15
## #18 0.04270852 74164.91  15
## #19 0.02630954 74223.75  15
```

```
# get the coefficients of the model with the best tau / lambda
coef(lasso_model, s = lasso_model$lambda.min)
```

```
## 16 x 1 sparse Matrix of class "dgCMatrix"
##              (Intercept) -5.072255e+03
## H              7.184295e+01
## So             4.466407e+01
## Ed             1.253075e+02
## Poi            1.023402e+02
## LF              .
## Po2             .
## H.F            6.888147e+01
## Pop            1.831089e+01
## NM             -2.143645e+03
## U1             8.835503e+01
## Wealt          7.718072e+03
## Ineq           4.882540e+01
## Prob          -3.688177e+03
## Time           .
```

```
# get predictions with best model
y_pred <- predict(lasso_model, s = lasso_model$lambda.min, newx=X)
```

```
# calculate R-squared
ss_total <- sum((y-mean(y))^2)
ss_residual <- sum((y-y_pred)^2)
r_squared <- 1 - (ss_residual/ss_total)
cat("R-squared:", r_squared)
```

```
## R-squared: 0.7743174
```

Using the crime data set, build a regression model using Elastic Net Regression

Comparatively, I created an elastic net model using `cv.glmnet`, which automatically scaled the data for me. The tau / lambda threshold that worked best for the model was 6.705118, which best minimized the MSE, and resulted in a regression with 14 non-zero variables. In other words, the variable selection process removed 1 predictor. The model that best minimized the MSE resulted in an R² of 0.791126.

```
# set seed for reproducibility
set.seed(1)
```

```
# predictors and response
X = as.matrix(crime_df[,16])
y = as.matrix(crime_df[,16])

# do k-fold cross validation for Elastic Net Regression model
enet_model <- cv.glmnet(x = X
                       , y = y
                       , alpha = .50 # elastic regression alpha between 0,1
                       , nfolds = 8 # number of folds
                       , nlambda = 20 # tau thresholds randomly generated
                       , type.measure = "mse" # squared error for gaussian models
                       , family = "gaussian"
                       , standardize = TRUE # use automatically scaled data
                       )

# plot MSE of elastic net model
plot(enet_model)
```



```
# use the tau / lambda that corresponds to the lowest MSE
enet_model$lambda.min
```

```
## [1] 6.705118
```

```
# get a list of tau/lambda, cross-validation error, and number of non-zero coefficients for each lambda.
cbind(enet_model$lambda, enet_model$cvm, enet_model$zero)
```

```
##              [,1]      [,2] [,3]
## #0 526.19079328 151536.04  0
## #1 324.05365872 126536.67  2
## #2 199.56786601 108766.05  2
## #3 122.90351327 98435.70  4
## #4 75.68990878 83147.98  7
## #5 46.61344992 75893.91  11
## #6 28.79683737 72929.82  12
## #7 17.70850449 70850.82  13
## #8 10.88761408 68961.93  14
## #9 6.70511766 68416.98  12
## #10 4.1293471 69451.08  13
## #11 2.54304340 70813.59  13
## #12 1.56612873 72393.45  15
## #13 0.96449757 73240.82  15
## #14 0.59399411 73496.26  15
## #15 0.36580405 73670.16  15
## #16 0.22527977 73855.96  15
## #17 0.13873814 74000.55  14
## #18 0.08544164 74102.59  15
## #19 0.05261908 74168.14  15
```

```
# get the coefficients of the model with the best tau / lambda
coef(enet_model, s = enet_model$lambda.min)
```

```
## 16 x 1 sparse Matrix of class "dgCMatrix"
##              (Intercept) -5.813274e+03
## H              9.255858e+01
## So             3.983590e+01
## Ed             1.477661e+02
## Poi            9.635595e+01
## Po2             .
## LF              .
## H.F            1.967116e+01
## Pop            -3.459667e-01
## NM             1.833806e+00
## U1             -3.841129e+03
## U2            1.312960e+02
## Wealt          5.196518e-02
## Ineq           5.778515e+01
## Prob          -3.961445e+03
## Time           .
```

```
# get predictions with best model
y_pred2 <- predict(enet_model, s = enet_model$lambda.min, newx=X)
```

```
# calculate R-squared
ss_total2 <- sum((y-mean(y))^2)
ss_residual2 <- sum((y-y_pred2)^2)
r_squared2 <- 1 - (ss_residual2/ss_total2)
cat("R-squared:", r_squared2)
```

```
## R-squared: 0.791126
```

It is interesting that all three methods removed the variable `Time`.

Out of the two global optimization variable selection methods, the Lasso model resulted in a model with a lower error at 67598.55 while the Elastic Net model had a slightly higher error of 68930.85, but the Elastic Net model had a higher R² value of 0.791126 versus the 0.7743174 of Lasso.