

Homework 5 - Coding

John L Kaspers
June 10, 2020

Warning

Warning: If you get an error when running the code, make sure that you have the necessary packages installed. You can install these packages by running the following code in base R (not RStudio or RMarkdown).

```
install.packages("air4")
install.packages("cvTools")
install.packages("glmnet")
install.packages("sandwich")
```

Model Comparison (5 points)

Use the `Puromycin` data set built into R. Further information on the data set can be found [here](#).

For this problem, you will run through various methods of model comparison - Adjusted R-squared, AIC, LRT, and LOOCV.

- Create two models that predict `rate`. The first model will use `conc` as the predictor. The second model will use `conc` and `state` as predictors. Do not include any interactions of other transformations. Run each of your models through the `summary()` function.

```
data(Puromycin)

lm.predict.rate1 <- lm(rate ~ conc, data = Puromycin)

lm.predict.rate2 <- lm(rate ~ conc + state, data = Puromycin)

summary(lm.predict.rate1)

##
## Call:
## lm(formula = rate ~ conc, data = Puromycin)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -49.861 -15.247  -2.861  15.686  48.054
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    93.92      8.00   11.74 1.09e-10 ***
## conc          105.48     16.92   6.23 3.53e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 28.82 on 21 degrees of freedom
## Multiple R-squared:  0.6489, Adjusted R-squared:  0.6322
## F-statistic: 38.81 on 1 and 21 DF,  p-value: 3.526e-06

summary(lm.predict.rate2)

##
## Call:
## lm(formula = rate ~ conc + state, data = Puromycin)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -61.381 -16.502   4.268  21.346  37.452
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    106.138      9.413  11.296 3.92e-10 ***
## conc          102.160     15.721   6.498 2.48e-06 ***
## stateuntreated -23.844     11.177  -2.133  0.0455 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26.65 on 20 degrees of freedom
## Multiple R-squared:  0.714, Adjusted R-squared:  0.6854
## F-statistic: 24.96 on 2 and 20 DF,  p-value: 3.667e-06
```

Which model is preferred according to adjusted R-squared?

According to adjusted R-squared, the model with `conc` and `state` as predictors is preferred because it has a slightly larger adjusted R-squared (0.6854 compared to 0.6322).

- Use the `AIC()` function to calculate the AIC for each model.

```
AIC(lm.predict.rate1)

## [1] 223.7834

AIC(lm.predict.rate2)

## [1] 221.0681
```

Which model is preferred according to AIC?

Because a lower AIC represents a better fitting model, model 2 (the model with `conc` and `state` as predictors) is preferred according to AIC because model 2 has a lower AIC (221.0681 compared to model 1's 223.7834).

- Run a likelihood ratio test between the two models.

```
anova(lm.predict.rate1, lm.predict.rate2, test = "LRT")

## Analysis of Variance Table
##
## Model 1: rate ~ conc
## Model 2: rate ~ conc + state
## Res.Df  RSS Df Sum of Sq Pr(>Chi)
## 1      21 17439
## 2      20 14206  1    3232.5  0.0329 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Which model is preferred according to the LRT?

According to LRT, the results are statistically significant as shown through the `p-value` of 0.0329, which is statistically significant at the 5% level so we should improve model 1 by adding 'state', thus model 2 is preferred according to the LRT.

- Use the `cvFit` function within the `cvTools` package to run LOOCV for both models.

```
cvFit(lm.predict.rate1, data = Puromycin, y = Puromycin$rate, cost = rmse, K = nrow(Puromycin))

## Leave-one-out CV results:
##      CV
## 30.66056

cvFit(lm.predict.rate2, data = Puromycin, y = Puromycin$rate, cost = rmse, K = nrow(Puromycin))

## Leave-one-out CV results:
##      CV
## 29.01072
```

Which model is preferred according to LOOCV?

According to LOOCV, in which we want a lower RMSPE for a better fitting model, `lm.predict.rate2` is marginally preferred because it has a slightly lower RMSPE (29.01072 compared to model 1's 30.66056).

Lasso (5 points)

Use the `swiss` data set built into R. Further information on the data set can be found [here](#).

For this problem, you will use Lasso regression to predict `Fertility` based on all other predictors. You will need to utilize the `glmnet()` and `cv.glmnet()` functions within the `glmnet` package. See the Lecture Notes for details.

- Create the design matrix using the `model.matrix()` function. This is used as one of the arguments for `glmnet()` when running Lasso.

```
data(swiss)
model.matrix(Fertility ~ ., data = swiss)

##              (Intercept) Agriculture Examination Education Catholic
## Courtelary             1      17.0      15      12      9.96
## Delemont                1      45.1       6      19      84.84
## Franches-Mnt            1      39.7       5       5      93.40
## Moutier                 1      36.5      12       7      33.77
## Neucheville             1      43.5      17      15      51.16
## Porrentruy              1      35.3       9       7      90.57
## Broye                   1      70.2      16       7      92.85
## Glane                   1      67.8      14       8      97.16
## Gruyere                 1      53.3      12       7      97.67
## Sarine                  1      45.2      16      13      91.28
## Veveyse                 1      64.5      14       6      98.61
## Aigle                   1      62.0      21      12      8.52
## Aubonne                 1      67.5      14       7      3.30
## Avenches                1      60.7      19      12      4.43
## Cossonay                1      69.3      22       5      2.82
## Echallens               1      72.6      18       2      24.20
## Grandson                1      34.0      17       8      3.30
## Lausanne                1      19.4      26      28     12.11
## La Vallee               1      15.2      31      20      2.15
## Lavaux                  1      73.0      19       9      2.84
## Morges                  1      59.8      22      10      5.23
## Noudon                  1      55.1      14       3      4.52
## Nyone                   1      50.9      22      12     15.14
## Orbe                    1      54.1      20       6      4.20
## Oron                    1      71.2      12      11      2.40
## Payerne                 1      58.1      14       8      5.23
## Paysd'enhaut            1      63.5       6       3      2.56
## Rolle                   1      60.8      16      10      7.72
## Vevey                   1      26.8      25      19     18.46
## Yverdon                 1      49.5      15       8      6.10
## Conthey                 1      85.9       3       2     99.71
## Entremont               1      84.9       7       6     99.68
## Herens                  1      89.7       5      22     100.00
## Martigny                1      78.2      12       6     98.96
## Monthey                 1      64.9       7      38     98.22
## St Maurice              1      75.9       9       9     99.06
## Sierre                  1      84.6       3       3     99.46
## Sion                    1      63.1      13      13     96.83
## Boudry                  1      38.4      26      12      5.62
## La Chauxd'fnd           1      7.7       29     11     13.79
## Le Locle                 1      16.7      22      13     11.22
## Neuchatel               1      17.6      35      32     16.92
## Val de Ruz               1      37.6      15       7      4.97
## ValdeTraversa           1      18.7      25      17      8.65
## V. De Geneve            1       1.2      37      53     42.34
## Rive Droite              1      46.6      16      29     50.43
## Rive Gauche             1      27.7      22      29     58.33
##
## Infant.Mortality
## Courtelary             22.2
## Delemont               22.2
## Franches-Mnt           20.2
## Moutier                20.3
## Neucheville            20.6
## Porrentruy             26.6
## Broye                  23.6
## Glane                   24.9
## Gruyere                 21.0
## Sarine                  24.4
## Veveyse                 24.5
## Aigle                   16.5
## Aubonne                 19.1
## Avenches                22.7
## Cossonay                18.7
## Echallens              21.2
## Grandson                20.0
## Lausanne                20.2
## La Vallee               10.8
## Lavaux                  20.0
## Morges                  19.0
## Noudon                  22.4
## Nyone                   16.7
## Orbe                    15.3
## Oron                    21.0
## Payerne                 23.8
## Paysd'enhaut            18.0
## Rolle                   16.3
## Vevey                   20.8
## Yverdon                 22.5
## Conthey                 15.1
## Entremont              19.8
## Herens                  18.3
## Martigny               19.4
## Monthey                 20.2
## St Maurice              17.8
## Sierre                  16.3
## Sion                    18.1
## Boudry                  20.3
## La Chauxd'fnd           20.5
## Le Locle                 18.9
## Neuchatel               23.0
## Val de Ruz              20.0
## ValdeTraversa           19.5
## V. De Geneve            18.0
## Rive Droite             18.2
## Rive Gauche             19.3
## attr(,"assign")
## [1] 0 1 2 3 4 5
```

- Run Lasso using different values of λ to determine which variable is the first to have its estimated slope coefficient go to 0. This may take multiple iterations. Remember to set the values of α to 1, as we did in the Lecture Notes.

```
library(glmnet)

lasso.model <- glmnet(model.matrix(Fertility ~ ., data = swiss), swiss$Fertility, alpha = 1)
coef(lasso.model, c(0, 0.1, 1, 1.05, 1.06, 10, 60, 100))

## 7 x 8 sparse Matrix of class "dgCMatrix"
##              1      2      3      4
## (Intercept)  66.6423597 65.7665041 55.8153858040 56.0341537090
## (Intercept) -0.1678682 -0.1544567 -0.0007866263 -0.0001253665
## Agriculture -0.2558537 -0.2467932 -0.1407473777 -0.1417288911
## Examination -0.8640181 -0.8430649 -0.6021498969 -0.5964469614
## Education    0.1031075  0.1000933  0.0653707291  0.0642655815
## Catholic     1.0764301  1.0734722  1.0334226904  1.0207250455
## Infant.Mortality 5      6      7      8
## (Intercept)  56.07860107 70.14255 70.14255 70.14255
## (Intercept)  .      .      .      .
## Agriculture  .      .      .      .
## Examination -0.14193272 .      .      .
## Education   -0.59531230 .      .      .
## Catholic    0.06404512 .      .      .
## Infant.Mortality 1.01817652 .      .      .
```

Which variable was the first to have its estimated slope coefficient go to 0?

Agriculture was the first variable to have its estimated slope coefficient go to 0.

- At what (approximate) value of λ do all of the estimated slope coefficients go to 0? An integer value is fine here!

```
cv.glmnet(model.matrix(Fertility ~ ., data = swiss), swiss$Fertility, alpha = 1)

##
## Call: cv.glmnet(x = model.matrix(Fertility ~ ., data = swiss), y = swiss$Fertility, alpha = 1)
##
## Measure: Mean-Squared Error
##
## Lambda Measure      SE Nonzero
## min 0.0256  59.66  9.30  5
## 1se 1.4006  68.58 11.16  4

60

4. Use the cv.glmnet() function to determine the optimal value for  $\lambda$ .
```

```
cv.glmnet(model.matrix(Fertility ~ ., data = swiss), swiss$Fertility, alpha = 1)$lambda.min

## [1] 0.04917668
```

What value of λ is returned? Note: This will vary from trial-to-trial (that's okay).

- Run lasso regression again using your optimal λ from Part 4.

```
coef(lasso.model, 0.04917668)

## 7 x 1 sparse Matrix of class "dgCMatrix"
##              1
## (Intercept)  66.3499116
## (Intercept)  .
## Agriculture  -0.1633465
## Examination -0.2530268
## Education    -0.8568413
## Catholic     0.1020699
## Infant.Mortality 1.0754606
```

What variables are in the model?

Agriculture, Examination, Education, Catholic and Infant.Mortality.

Variances (12 points)

Use the stopping data set from the `air4` package. Further information on the data set can be found [here](#).

For this problem, you will run through various methods of dealing with nonconstant variance - using OLS, WLS (assuming known weights), WLS (with a sandwich estimator), and bootstrap.

- Create a scatterplot of Distance versus Speed (Y vs X).

```
data(stopping, package = "air4")

plot(stopping$Speed, stopping$Distance,
     main = "Distance vs. Speed",
     xlab = "Speed",
     ylab = "Distance")

##
## Distance vs. Speed
##
##      Speed      Distance
## 1      5         10
## 2      5         10
## 3      5         10
## 4      5         10
## 5      5         10
## 6      5         10
## 7      5         10
## 8      5         10
## 9      5         10
## 10     5         10
## 11     5         10
## 12     5         10
## 13     5         10
## 14     5         10
## 15     5         10
## 16     5         10
## 17     5         10
## 18     5         10
## 19     5         10
## 20     5         10
## 21     5         10
## 22     5         10
## 23     5         10
## 24     5         10
## 25     5         10
## 26     5         10
## 27     5         10
## 28     5         10
## 29     5         10
## 30     5         10
## 31     5         10
## 32     5         10
## 33     5         10
## 34     5         10
## 35     5         10
## 36     5         10
## 37     5         10
## 38     5         10
## 39     5         10
## 40     5         10
## 41     5         10
## 42     5         10
## 43     5         10
## 44     5         10
## 45     5         10
## 46     5         10
## 47     5         10
## 48     5         10
## 49     5         10
## 50     5         10
## 51     5         10
## 52     5         10
## 53     5         10
## 54     5         10
## 55     5         10
## 56     5         10
## 57     5         10
## 58     5         10
## 59     5         10
## 60     5         10
## 61     5         10
## 62     5         10
## 63     5         10
## 64     5         10
## 65     5         10
## 66     5         10
## 67     5         10
## 68     5         10
## 69     5         10
## 70     5         10
## 71     5         10
## 72     5         10
## 73     5         10
## 74     5         10
## 75     5         10
## 76     5         10
## 77     5         10
## 78     5         10
## 79     5         10
## 80     5         10
## 81     5         10
## 82     5         10
## 83     5         10
## 84     5         10
## 85     5         10
## 86     5         10
## 87     5         10
## 88     5         10
## 89     5         10
## 90     5         10
## 91     5         10
## 92     5         10
## 93     5         10
## 94     5         10
## 95     5         10
## 96     5         10
## 97     5         10
## 98     5         10
## 99     5         10
## 100    5         10
## 101    5         10
## 102    5         10
## 103    5         10
## 104    5         10
## 105    5         10
## 106    5         10
## 107    5         10
## 108    5         10
## 109    5         10
## 110    5         10
## 111    5         10
## 112    5         10
## 113    5         10
## 114    5         10
## 115    5         10
## 116    5         10
## 117    5         10
## 118    5         10
## 119    5         10
## 120    5         10
## 121    5         10
## 122    5         10
## 123    5         10
## 124    5         10
## 125    5         10
## 126    5         10
## 127    5         10
## 128    5         10
## 129    5         10
## 130    5         10
## 131    5         10
## 132    5         10
## 133    5         10
## 134    5         10
## 135    5         10
## 136    5         10
## 137    5         10
## 138    5         10
## 139    5         10
## 140    5         10
## 141    5         10
## 142    5         10
## 143    5         10
## 144    5         10
## 145    5         10
## 146    5         10
## 147    5         10
## 148    5         10
## 149    5         10
## 150    5         10
## 151    5         10
## 152    5         10
## 153    5         10
## 154    5         10
## 155    5         10
## 156    5         10
## 157    5         10
## 158    5         10
## 159    5         10
## 160    5         10
## 161    5         10
## 162    5         10
## 163    5         10
## 164    5         10
## 165    5         10
## 166    5         10
## 167    5         10
## 168    5         10
## 169    5         10
## 170    5         10
## 171    5         10
## 172    5         10
## 173    5         10
## 174    5         10
## 175    5         10
## 176    5         10
## 177    5         10
## 178    5         10
## 179    5         10
## 180    5         10
## 181    5         10
## 182    5         10
## 183    5         10
## 184    5         10
## 185    5         10
## 186    5         10
## 187    5         10
## 188    5         10
## 189    5         10
## 190    5         10
## 191    5         10
## 192    5         10
## 193    5         10
## 194    5         10
## 195    5         10
## 196    5         10
## 197    5         10
## 198    5         10
## 199    5         10
## 200    5         10
## 201    5         10
## 202    5         10
## 203    5         10
## 204    5         10
## 205    5         10
## 206    5         10
## 207    5         10
## 208    5         10
## 209    5         10
## 210    5         10
## 211    5         10
## 212    5         10
## 213    5         10
## 214    5         10
## 215    5         10
## 216    5         10
## 217    5         10
## 218    5         10
## 219    5         10
## 220    5         10
## 221    5         10
## 222    5         10
## 223    5         10
## 224    5         10
## 225    5         10
## 226    5         10
## 227    5         10
## 228    5         10
## 229    5         10
## 230    5         10
## 231    5         10
## 232    5         10
## 233    5         10
## 234    5         10
## 235    5         10
## 236    5         10
## 237    5         10
## 238    5         10
## 239    5         10
## 240    5         10
## 241    5         10
## 242    5         10
## 243    5         10
## 244    5         10
## 245    5         10
## 246    5         10
## 247    5         10
## 248    5         10
## 249    5         10
## 250    5         10
## 251    5         10
## 252    5         10
## 253    5         10
## 254    5         10
## 255    5         10
## 256    5         10
## 257    5         10
## 258    5         10
## 259    5         10
## 260    5         10
## 261    5         10
## 262    5         10
## 263    5         10
## 264    5         10
## 265    5         10
## 266    5         10
## 267    5         10
## 268    5         10
## 269    5         10
## 270    5         10
## 271    5         10
## 272    5         10
## 273    5         10
## 274    5         10
## 275    5         10
## 276    5         10
## 277    5         10
## 278    5         10
## 279    5         10
## 280    5         10
## 281    5         10
## 282    5         10
## 283    5         10
## 284    5         10
## 285    5         10
## 286    5         10
## 287    5         10
## 288    5         10
## 289    5         10
## 290    5         10
## 291    5         10
## 292    5         10
## 293    5         10
## 294    5         10
## 295    5         10
## 296    5         10
## 297    5         10
## 298    5         10
## 299    5         10
## 300    5         10
## 301    5         10
## 302    5         10
## 303    5         10
## 304    5         10
## 305    5         10
## 306    5         10
## 307    5         10
## 308    5         10
## 309    5         10
## 310    5         10
## 311    5         10
## 312    5         10
## 313    5         10
## 314    5         10
## 315    5         10
## 316    5         10
## 317    5         10
## 318    5         10
## 319    5         10
## 320    5         10
## 321    5         10
## 322    5         10
## 323    5         10
## 324    5         10
## 325    5         10
## 326    5         10
## 327    5         10
## 328    5         10
## 329    5         10
## 330    5         10
## 331    5         10
## 332    5         10
## 333    5         10
## 334    5         10
## 335    5         10
## 336    5         10
## 337    5         10
## 338    5         10
## 339    5         10
## 340    5         10
## 341    5         10
## 342    5         10
## 343    5         10
## 344    5         10
## 345    5         10
## 346    5         10
## 347    5         10
## 348    5         10
## 349    5         10
## 350    5         10
## 351    5         10
## 352    5         10
## 353    5         10
## 354    5         10
## 355    5         10
## 356    5         10
## 357    5         10
## 358    5         10
## 359    5         10
## 360    5         10
## 361    5         10
## 362    5         10
## 363    5         10
## 364    5         10
## 365    5         10
## 366    5         10
## 367    5         10
## 368    5         10
## 369    5         10
## 370    5         10
## 371    5         10
## 372    5         10
## 373    5         10
## 374    5         10
## 375    5         10
## 376    5         10
## 377    5         10
## 378    5         10
## 379    5         10
## 380    5         10
## 381    5         10
## 382    5         10
## 383    5         10
## 384    5         10
## 385    5         10
## 386    5         10
## 387    5         10
## 388    5         10
## 389    5         10
## 390    5         10
## 391    5         10
## 392    5         10
## 393    5         10
## 394    5         10
## 395    5         10
## 396    5         10
## 397    5         10
## 398    5         10
## 399    5         10
## 400    5         10
## 401    5         10
## 402    5         10
## 403    5         10
## 404    5         10
## 405    5         10
## 406    5         10
## 407    5         10
## 408    5         10
## 409    5         10
## 410    5         10
## 411    5         10
## 412    5         10
## 413    5         10
## 414    5         10
## 415    5         10
## 416    5         10
## 417    5         10
## 418    5         10
## 419    5         10
## 420    5         10
## 421    5         10
## 422    5         10
## 423    5         10
## 424    5         10
## 425    5         10
## 426    5         10
## 427    5         10
## 428    5         10
## 429    5         10
## 430    5         10
## 431    5         10
## 432    5         10
## 433    5         10
## 434    5         10
## 435    5         10
## 436    5         10
## 437    5         10
## 438    5         10
## 439    5         10
## 440    5         10
## 441    5         10
## 442    5         10
## 443    5         10
## 444    5         10
## 445    5         10
## 446    5         10
## 447    5         10
## 448    5         10
## 449    5         10
## 450    5         10
## 451    5         10
## 452    5         10
## 453    5         10
## 454    5         10
## 455    5         10
## 456    5         10
## 457    5         10
## 458    5         10
## 459    5         10
## 460    5         10
## 461    5         10
## 462    5         10
## 463    5         10
## 464    5         10
## 465    5         10
## 466    5         10
## 467    5         10
## 468    5         10
## 469    5         10
## 470    5         10
## 471    5         10
## 472    5         10
## 473    5         10
## 474    5         10
## 475    5         10
```