```
1   > set.seed(1)
2   >
3   > #install.packages("stats")
4   > #install.packages("DAAG")
5   > library(stats)
6   > library(DAAG)
7   >
8   > input = data.frame(read.table("uscrime.txt", header = TRUE)) #read in data
9   > mydata = input[c(16, 1:15)] #reorder so that crime is the first column (for formula)
10  >
11  > # Plot predictors vs. response
12  > predictors = mydata[-1]
13  > headers = list(
14  +    "M",
15  +    "So",
16  +    "Ed",
17  +    "Po1",
18  +    "Po2",
19  +    "LF",
20  +    "M.F",
21  +    "Pop",
22  +    "NW",
23  +    "U1",
24  +    "U2",
25  +    "Wealth",
26  +    "Ineq",
27  +    "Prob",
28  +    "Time"
29  + )
30  > par(mfrow = c(4, 4))
31  > for (i in 1:15) {
32  +    plot(predictors[, i], mydata$Crime, xlab = headers[i])
33  + }
34  >
35  >
36  > point = data.frame(
37  +    M = 14.0,
38  +    So = 0,
39  +    Ed = 10.0,
40  +    Po1 = 12.0,
41  +    Po2 = 15.5,
42  +    LF = 0.640,
43  +    M.F = 94.0,
44  +    Pop = 150,
45  +    NW = 1.1,
46  +    U1 = 0.120,
47  +    U2 = 3.6,
48  +    Wealth = 3200,
49  +    Ineq = 20.1,
50  +    Prob = 0.04,
51  +    Time = 39.0
52  + )
53  >
54  >
55  > f1 = formula(mydata)
56  > model1 = lm(f1, mydata)
57  > summary(model1)
58
59  Call:
60  lm(formula = f1, data = mydata)
61
62  Residuals:
63      Min      1Q Median      3Q     Max
64  -395.7   -98.1   -6.7   113.0   512.7
65
66  Coefficients:
```

```
67                   Estimate Std. Error t value Pr(>|t|)
68   (Intercept) -5.98e+03   1.63e+03    -3.68  0.00089 ***
69   M             8.78e+01   4.17e+01     2.11  0.04344 *
70   So           -3.80e+00   1.49e+02    -0.03  0.97977
71   Ed            1.88e+02   6.21e+01     3.03  0.00486 **
72   Po1           1.93e+02   1.06e+02     1.82  0.07889 .
73   Po2          -1.09e+02   1.17e+02    -0.93  0.35883
74   LF           -6.64e+02   1.47e+03    -0.45  0.65465
75   M.F           1.74e+01   2.04e+01     0.86  0.39900
76   Pop          -7.33e-01   1.29e+00    -0.57  0.57385
77   NW            4.20e+00   6.48e+00     0.65  0.52128
78   U1           -5.83e+03   4.21e+03    -1.38  0.17624
79   U2            1.68e+02   8.23e+01     2.04  0.05016 .
80   Wealth        9.62e-02   1.04e-01     0.93  0.36075
81   Ineq          7.07e+01   2.27e+01     3.11  0.00398 **
82   Prob         -4.86e+03   2.27e+03    -2.14  0.04063 *
83   Time         -3.48e+00   7.17e+00    -0.49  0.63071
84   ---
85   Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
86
87   Residual standard error: 209 on 31 degrees of freedom
88   Multiple R-squared:  0.803, Adjusted R-squared:  0.708
89   F-statistic: 8.43 on 15 and 31 DF,  p-value: 3.54e-07
90
91   > coef1 = model1$coefficients
92   >
93   > par(mfrow = c(2, 2))
94   > plot(model1)
95   >
96   > crime_prediction = predict.lm(model1, point)
97   > crime_prediction
98     1
99   155
100  > #that answer of 155 is really low, we are probably overfit, so let's look at the
     p-values of each point.
101  >
102  >
103  > Pvalues = summary(model1)$coefficients[, 4]
104  > coef = model1$coefficients
105  >
106  >
107  > # Eliminate those predictors with a p-value > 0.08.  I know  0.05 is usually the rule,
108  > # but the U2 factor (unemployment rate of urban males 35-39) and Po1
109  > # should be left in as it was considered important by the summary() function.
110  > mydata_fit = mydata[1]
111  > n = 2
112  > for (i in 2:16) {
113  +   if (Pvalues[i] < 0.08) {
114  +     mydata_fit[n] = mydata[i]
115  +     n = n + 1
116  +   }
117  + }
118  >
119  > f2 = formula(mydata_fit)
120  > model2 = lm(f2, mydata_fit)
121  > summary(model2)
122
123  Call:
124  lm(formula = f2, data = mydata_fit)
125
126  Residuals:
127      Min    1Q Median    3Q    Max
128  -470.7  -78.4  -19.7  133.1  556.2
129
130  Coefficients:
131                Estimate Std. Error t value Pr(>|t|)
```

```
132    (Intercept)  -5040.5       899.8   -5.60  1.7e-06 ***
133    M              105.0        33.3    3.15   0.0031 **
134    Ed             196.5        44.8    4.39  8.1e-05 ***
135    Po1            115.0        13.8    8.36  2.6e-10 ***
136    U2              89.4        40.9    2.18   0.0348 *
137    Ineq            67.7        13.9    4.85  1.9e-05 ***
138    Prob         -3801.8      1528.1   -2.49   0.0171 *
139    ---
140    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
141
142    Residual standard error: 201 on 40 degrees of freedom
143    Multiple R-squared:  0.766, Adjusted R-squared:  0.731
144    F-statistic: 21.8 on 6 and 40 DF,  p-value: 3.42e-11
145
146    > plot(model2)
147    >
148    >
149    > crime_prediction_adj = predict.lm(model2, point)
150    >
151    > crime_prediction_adj
152       1
153    1304
154    >
155    > #Try cross-validation as well:
156    > par(mfrow = c(1, 1))
157    > c1 = cv.lm(mydata, model1, m = 5)
158    Analysis of Variance Table
159
160    Response: Crime
161              Df  Sum Sq Mean Sq F value   Pr(>F)
162    M          1    55084   55084    1.26  0.2702
163    So         1    15370   15370    0.35  0.5575
164    Ed         1   905668  905668   20.72 7.7e-05 ***
165    Po1        1  3076033 3076033   70.38 1.8e-09 ***
166    Po2        1   153024  153024    3.50  0.0708 .
167    LF         1    61134   61134    1.40  0.2459
168    M.F        1   111000  111000    2.54  0.1212
169    Pop        1    42649   42649    0.98  0.3309
170    NW         1    14197   14197    0.32  0.5728
171    U1         1     7065    7065    0.16  0.6904
172    U2         1   269663  269663    6.17  0.0186 *
173    Wealth     1    34748   34748    0.79  0.3795
174    Ineq       1   547423  547423   12.52  0.0013 **
175    Prob       1   222620  222620    5.09  0.0312 *
176    Time       1    10304   10304    0.24  0.6307
177    Residuals 31  1354946   43708
178    ---
179    Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
180
181
182
183    fold 1
184    Observations in test set: 9
185                    1    4    8    9   18      20    23    32    47
186    Predicted     755 1791 1362 689  844 1227.84  958 807.8  992
187    cvpred        658 1690 1300 617  792 1220.22  814 804.9 1077
188    Crime         791 1969 1555 856  929 1225.00 1216 754.0  849
189    CV residual   133  279  255 239  137    4.78  402 -50.9 -228
190
191    Sum of squares = 453204    Mean square = 50356    n = 9
192
193    fold 2
194    Observations in test set: 10
195                    5   13   15   17   25    34    39     40    42   46
196    Predicted    1167  733  903  393  606 971.5 839.3 1131.5 326.3  827
197    cvpred       1132  926  977  152  740 902.7 918.1 1248.5  62.3 1004
```

```
198   Crime         1234  511  798 539  523 923.0 826.0 1151.0 542.0  508
199   CV residual   102 -415 -179 387 -217  20.3 -92.1  -97.5 479.7 -496
200
201   Sum of squares = 906384    Mean square = 90638    n = 10
202
203   fold 3
204   Observations in test set: 10
205                    2    3   11   14   16   22   28    31   33     38
206   Predicted   1473.7 322 1161  780 1006  657 1258 388.0  841 562.693
207   cvpred      1566.9 313  953  782 1129  876 1368 321.7  700 566.231
208   Crime       1635.0 578 1674  664  946  439 1216 373.0 1072 566.000
209   CV residual   68.1 265  721 -118 -183 -437 -152  51.3  372  -0.231
210
211   Sum of squares = 997216    Mean square = 99722    n = 10
212
213   fold 4
214   Observations in test set: 9
215                   19    21   26   27   29   30    36   44   45
216   Predicted   1146 774.9 1977  279 1287 702.7 1137.6 1121  617
217   cvpred      1529 802.3 1673  467 1673 629.6 1191.9 1298  702
218   Crime        750 742.0 1993  342 1043 696.0 1272.0 1030  455
219   CV residual -779 -60.3  320 -125 -630  66.4   80.1 -268 -247
220
221   Sum of squares = 1269688    Mean square = 141076    n = 9
222
223   fold 5
224   Observations in test set: 9
225                    6    7   10   12   24   35   37   41   43
226   Predicted    793 934.2 736.5 722.0 869 737.8  971 824 1134
227   cvpred       819 950.9 758.1 772.5 802 690.5 1227 891 1267
228   Crime        682 963.0 705.0 849.0 968 653.0  831 880  823
229   CV residual -137  12.1 -53.1  76.5 166 -37.5 -396 -11 -444
230
231   Sum of squares = 410109    Mean square = 45568    n = 9
232
233   Overall (Sum over all 9 folds)
234      ms
235   85885
236   Warning message:
237   In cv.lm(mydata, model1, m = 5) :
238
239    As there is >1 explanatory variable, cross-validation
240    predicted values for a fold are not a linear function
241    of corresponding overall predicted values.  Lines that
242    are shown for the different folds are approximate
243
244   > c2 = cv.lm(mydata, model2, m = 5)
245   Analysis of Variance Table
246
247   Response: Crime
248             Df  Sum Sq Mean Sq F value  Pr(>F)
249   M          1   55084   55084    1.37 0.24914
250   Ed         1  725967  725967   18.02 0.00013 ***
251   Po1        1 3173852 3173852   78.80 5.3e-11 ***
252   U2         1  217386  217386    5.40 0.02534 *
253   Ineq       1  848273  848273   21.06 4.3e-05 ***
254   Prob       1  249308  249308    6.19 0.01711 *
255   Residuals 40 1611057   40276
256   ---
257   Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
258
259
260
261   fold 1
262   Observations in test set: 9
263                    1    4    8    9   18    20   23    32   47
```

```
264  Predicted    810.8 1897 1354 719 800 1203.0  938 773.7  976
265  cvpred       762.1 1858 1282 657 672 1210.8  871 777.6  998
266  Crime        791.0 1969 1555 856 929 1225.0 1216 754.0  849
267  CV residual  28.9  111  273 199 257   14.2  345 -23.6 -149
268
269  Sum of squares = 335463    Mean square = 37274    n = 9
270
271  fold 2
272  Observations in test set: 10
273                  5    13    15    17   25    34    39    40  42   46
274  Predicted    1270   739 828.34 527.4  579   998 786.7 1141 369  748
275  cvpred       1337   842 804.73 469.3  671  1032 810.3 1187 302  839
276  Crime        1234   511 798.00 539.0  523   923 826.0 1151 542  508
277  CV residual  -103  -331  -6.73  69.7 -148  -109  15.7  -36 240 -331
278
279  Sum of squares = 327423    Mean square = 32742    n = 10
280
281  fold 3
282  Observations in test set: 10
283                  2    3    11    14     16    22     28    31    33    38
284  Predicted    1388 386 1118 713.6 1004.4   728 1259.0 440.4  874 544.4
285  cvpred       1368 390 1019 711.8  985.8   767 1252.6 423.8  850 511.2
286  Crime        1635 578 1674 664.0  946.0   439 1216.0 373.0 1072 566.0
287  CV residual   267 188  655 -47.8  -39.8  -328  -36.6 -50.8  222  54.8
288
289  Sum of squares = 702726    Mean square = 70273    n = 10
290
291  fold 4
292  Observations in test set: 9
293                 19    21    26    27    29    30   36   44   45
294  Predicted    1221 783.3 1789.1 312.20 1495 668.0 1102 1178  622
295  cvpred       1316 836.4 1895.7 334.15 1693 631.2 1163 1191  612
296  Crime         750 742.0 1993.0 342.00 1043 696.0 1272 1030  455
297  CV residual  -566 -94.4   97.3   7.85 -650  64.8  109 -161 -157
298
299  Sum of squares = 827924    Mean square = 91992    n = 9
300
301  fold 5
302  Observations in test set: 9
303                  6   7    10   12    24    35    37    41   43
304  Predicted    730 733 787.3 673 919.4   808   992 796.4 1017
305  cvpred       707 694 776.8 660 879.7   777  1115 812.6 1091
306  Crime        682 963 705.0 849 968.0   653   831 880.0  823
307  CV residual  -25 269 -71.8 189  88.3  -124  -284  67.4 -268
308
309  Sum of squares = 294201    Mean square = 32689    n = 9
310
311  Overall (Sum over all 9 folds)
312     ms
313  52931
314  Warning message:
315  In cv.lm(mydata, model2, m = 5) :
316
317   As there is >1 explanatory variable, cross-validation
318   predicted values for a fold are not a linear function
319   of corresponding overall predicted values.  Lines that
320   are shown for the different folds are approximate
321
322  >
323  > # Now compare the models using the R^2 values. From the summaries printed earlier
324  > # we know Model 1's R2 was 0.803 and Model 2's R2 was .766.
325  >
326  > SStot = sum((mydata$Crime - mean(mydata$Crime)) ^ 2)
327  >
328  > SSc1 = attr(c1, "ms") * nrow(mydata)
329  > SSc2 = attr(c2, "ms") * nrow(mydata)
```

```
> 
> R2_cvm1 = 1 - SSc1 / SStot
> R2_cvm2 = 1 - SSc2 / SStot
> R2_cvm1
[1] 0.413
> R2_cvm2
[1] 0.638
> 
> # So we see that the first model was overfit to the data. While the R2 of
> # model 1 was initially higher than model 2 using all of the data, by using
> # 5 fold cross validation we see that model 2 has a better fit, though it is
> # still probably over-fit as we only have a small set of data. As expected, the R2
> # of model 2 using cross validation is lower than that of the whole data set.
```