

Application

Refer to the Air Quality data described previously, and the analyses we have done with Ozone as the response variable, and the five explanatory variables (including the two engineered features).

2. Add GAM on all variables to the 10-fold CV comparison that has been used for LASSO, Ridge, and other methods. Use the same folds for GAM that were used for the other methods.

```
39 get.folds = function(n, K) {  
40   ### Get the appropriate number of fold labels  
41   n.fold = ceiling(n / K) # Number of observations per fold (rounded up)  
42   fold.ids.raw = rep(1:K, times = n.fold) # Generate extra labels  
43   fold.ids = fold.ids.raw[1:n] # Keep only the correct number of labels  
44  
45   ### Shuffle the fold labels  
46   folds.rand = fold.ids[sample.int(n)]  
47  
48   return(folds.rand)  
49 }  
50  
51  
52 ### Number of folds  
53 K = 10  
54  
55 ### Construct folds  
56 n = nrow(data) # Sample size  
57 folds = get.folds(n, K)  
58  
59 ### Create a container for MSPEs. Let's include ordinary least-squares  
60 ### regression for reference  
61 all.models = c("LS", "Hybrid", "Ridge", "LASSO-Min", "LASSO-lse", "GAM")  
62 all.MSPEs = array(0, dim = c(K, length(all.models)))  
63 colnames(all.MSPEs) = all.models  
64 ### Begin cross-validation  
65 for(i in 1:K){  
66   ### Split data  
67   data.train = data[folds != i,]  
68   data.valid = data[folds == i,]  
69   n.train = nrow(data.train)  
70  
71   ### Get response vectors  
72   Y.train = data.train$Ozone  
73   Y.valid = data.valid$Ozone
```

```

74
75 # LS
76 fit.ls = lm(Ozone ~ ., data = data.train)
77 pred.ls = predict(fit.ls, newdata = data.valid)
78 MSPE.ls = get.MSPE(Y.valid, pred.ls)
79 all.MSPEs[i, "LS"] = MSPE.ls
80
81 #Hybrid Stepwise
82
83 fit.start = lm(Ozone ~ 1, data = data.train)
84 fit.end = lm(Ozone ~ ., data = data.train)
85
86 step.BIC = step(fit.start, list(upper = fit.end), k = log(n.train),
87                      trace = 0)
88
89 pred.step.BIC = predict(step.BIC, data.valid)
90
91 err.step.BIC = get.MSPE(Y.valid, pred.step.BIC)
92
93 all.MSPEs[i, "Hybrid"] = err.step.BIC
94
95
96
97 #ridge regression
98 lambda.vals = seq(from = 0, to = 100, by = 0.05)
99
100 fit.ridge = lm.ridge(Ozone ~ ., lambda = lambda.vals,
101                      data = data.train)
102
103 ind.min.GCV = which.min(fit.ridge$GCV)
104 lambda.min = lambda.vals[ind.min.GCV]
105
106 all.coefs.ridge = coef(fit.ridge)
107 coef.min = all.coefs.ridge[ind.min.GCV,]
108
109 matrix.valid.ridge = model.matrix(Ozone ~ ., data = data.valid)

```

```

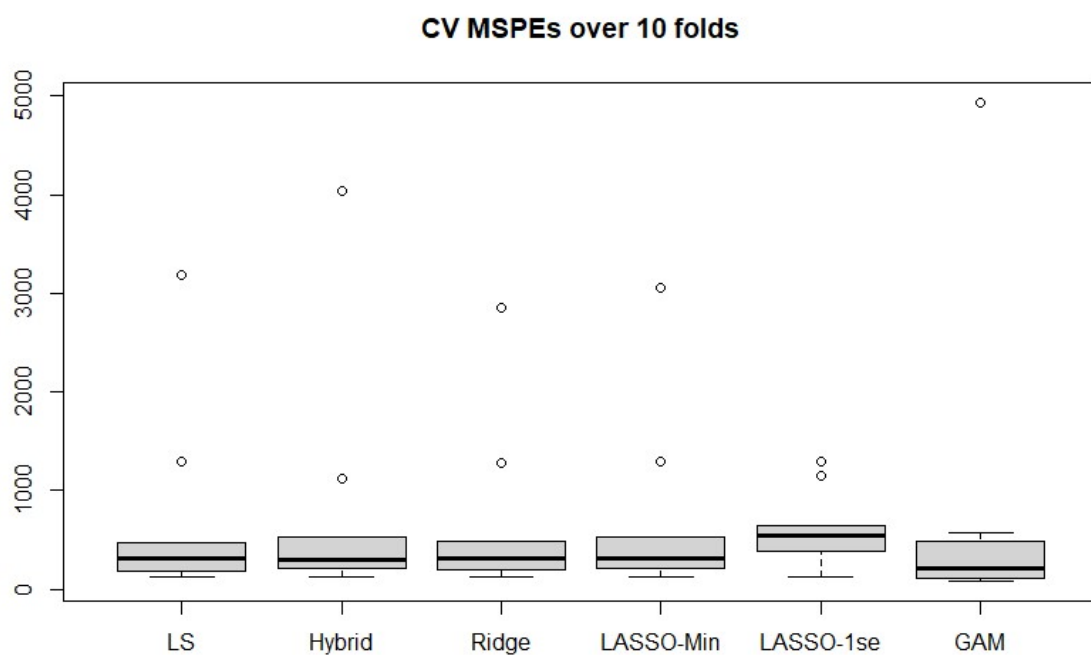
110
111 ### Now we can multiply the data by our coefficient vector. The
112 ### syntax in R for matrix-vector multiplication is %*%. Note that,
113 ### for this type of multiplication, order matters. That is,
114 ### A %*% B != B %*% A. Make sure you do data %*% coefficients.
115 ### For more information, see me in a Q&A session or, better still,
116 ### take a course on linear algebra (it's really neat stuff)
117 pred.ridge = matrix.valid.ridge %*% coef.min
118
119 ### Now we just need to calculate the MSPE and store it
120 MSPE.ridge = get.MSPE(Y.valid, pred.ridge)
121 all.MSPEs[i, "Ridge"] = MSPE.ridge
122
123 matrix.train.raw = model.matrix(Ozone ~ ., data = data.train)
124 matrix.train = matrix.train.raw[,-1]
125
126 ### LASSO
127 all.LASSOs = cv.glmnet(x = matrix.train, y = Y.train)
128
129 ### Get both 'best' lambda values using $lambda.min and $lambda.1se
130 lambda.min = all.LASSOs$lambda.min
131 lambda.1se = all.LASSOs$lambda.1se
132
133 ### Get the coefficients for our two 'best' LASSO models
134 coef.LASSO.min = predict(all.LASSOs, s = lambda.min, type = "coef")
135 coef.LASSO.1se = predict(all.LASSOs, s = lambda.1se, type = "coef")
136
137 ### Get which predictors are included in our models (i.e. which
138 ### predictors have non-zero coefficients)
139 included.LASSO.min = predict(all.LASSOs, s = lambda.min,
140                             type = "nonzero")
141 included.LASSO.1se = predict(all.LASSOs, s = lambda.1se,
142                             type = "nonzero")
143
144 matrix.valid.LASSO.raw = model.matrix(Ozone ~ ., data = data.valid)
145 matrix.valid.LASSO = matrix.valid.LASSO.raw[,-1]
146 pred.LASSO.min = predict(all.LASSOs, newx = matrix.valid.LASSO,
147                          s = lambda.min, type = "response")
148 pred.LASSO.1se = predict(all.LASSOs, newx = matrix.valid.LASSO,
149                          s = lambda.1se, type = "response")
150
151 ### Calculate MSPEs and store them
152 MSPE.LASSO.min = get.MSPE(Y.valid, pred.LASSO.min)
153 all.MSPEs[i, "LASSO-Min"] = MSPE.LASSO.min
154
155 MSPE.LASSO.1se = get.MSPE(Y.valid, pred.LASSO.1se)
156 all.MSPEs[i, "LASSO-1se"] = MSPE.LASSO.1se
157
158 ## GAM
159 fit.gam = gam(Ozone ~ s(Solar.R) + s(Wind) + s(Temp) + s(TWcp) + s(TWrat),
160              data = data.train)
161
162 pred.gam = predict(fit.gam, data.valid)
163 MSPE.gam = get.MSPE(Y.valid, pred.gam) # Our helper function
164 all.MSPEs[i, "GAM"] = MSPE.gam
165 }
166
167 all.MSPEs
168
169 ### Make a boxplot of MSPEs. I would like to include the number of folds
170 ### in the title. This can be done by using the paste0() function,
171 ### which concatenates strings (i.e. attaches them end-to-end), and
172 ### can be provided numeric variables.
173 boxplot(all.MSPEs, main = paste0("CV MSPEs over ", K, " folds"))

```

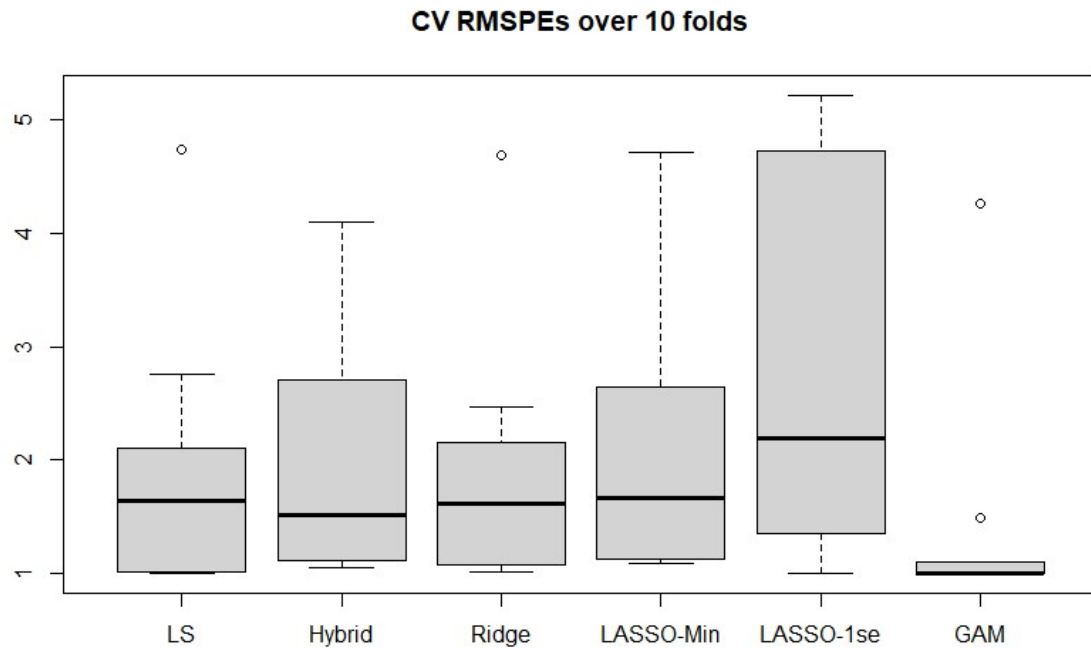
(a) Report the separate MSPEs from each fold, $MSPE_v$, $v = 1, \dots, 10$ and the MSPE for the full data.

GAM
 4938.2488
 80.6323
 192.6688
 571.8521
 207.6464
 274.0896
 80.5867
 207.3574
 488.7937
 109.5156

(b) Starting with boxplots the plots made earlier for least squares, hybrid stepwise, ridge, and LASSO, ADD a boxplot of the 10 CV error estimates for GAM as the last box on the right. Comment on how GAM compares to other methods



(c) Repeat this using relative MSPE.



(d) Using the knowledge gained from the analysis you did in Question 1, give a 1-sentence explanation for why GAM performs as it does. (If it is better than other methods, why? If it is no better than other methods, why?)

- ➔ Gam may not be the best model if the variables are interactive, however we include two interactions so it covers it.