

3. Add “tuned PPR” to the CV comparison that has been used for other methods. Use the same folds as were used for the other methods. Arrange it so that in each fold, the best PPR model is chosen exactly as above and is used to produce PPR’s predictions for that fold. *This means that you will need to tune PPR exactly as in the previous problem, separately within each fold!!! (You can use 5-fold CV for tuning if you want to save a little time. Separately save out the number of terms used in the best model for each fold.*

```
set.seed(50297026)
121 ### Number of folds
122 K = 10
123
124 ### Construct folds
125 n = nrow(data) # Sample size
126 folds = get.folds(n, K)
127
128 ### Create a container for MSPEs. Let's include ordinary least-squares
129 ### regression for reference
130 all.models = c("LS", "Hybrid", "Ridge", "LASSO-Min", "LASSO-lse", "GAM", "PPR")
131 all.MSPEs = array(0, dim = c(K, length(all.models)))
132 colnames(all.MSPEs) = all.models
133 ### Begin cross-validation
134 for(i in 1:K){
135   print(paste0(i, " of ", K))
136   ### Split data
137   data.train = data[folds != i,]
138   data.valid = data[folds == i,]
139   n.train = nrow(data.train)
140
141   ### Get response vectors
142   Y.train = data.train$Ozone
143   Y.valid = data.valid$Ozone
144
145   # LS
146   fit.ls = lm(Ozone ~ ., data = data.train)
147   pred.ls = predict(fit.ls, newdata = data.valid)
148   MSPE.ls = get.MSPE(Y.valid, pred.ls)
149   all.MSPEs[i, "LS"] = MSPE.ls
150
151   #Hybrid Stepwise
152
153   fit.start = lm(Ozone ~ 1, data = data.train)
154   fit.end = lm(Ozone ~ ., data = data.train)
155
156   step.BIC = step(fit.start, list(upper = fit.end), k = log(n.train),
157                  trace = 0)
158
159   pred.step.BIC = predict(step.BIC, data.valid)
160
161   err.step.BIC = get.MSPE(Y.valid, pred.step.BIC)
```

```

162
163 all.MSPEs[i, "Hybrid"] = err.step.BIC
164
165
166
167 #ridge regression
168 lambda.vals = seq(from = 0, to = 100, by = 0.05)
169
170 fit.ridge = lm.ridge(Ozone ~ ., lambda = lambda.vals,
171                      data = data.train)
172
173 ind.min.GCV = which.min(fit.ridge$GCV)
174 lambda.min = lambda.vals[ind.min.GCV]
175
176 all.coefs.ridge = coef(fit.ridge)
177 coef.min = all.coefs.ridge[ind.min.GCV,]
178
179 matrix.valid.ridge = model.matrix(Ozone ~ ., data = data.valid)
180
181 ### Now we can multiply the data by our coefficient vector. The
182 ### syntax in R for matrix-vector multiplication is %%. Note that,
183 ### for this type of multiplication, order matters. That is,
184 ### A %% B != B %% A. Make sure you do data %% coefficients.
185 ### For more information, see me in a Q&A session or, better still,
186 ### take a course on linear algebra (it's really neat stuff)
187 pred.ridge = matrix.valid.ridge %% coef.min
188
189 ### Now we just need to calculate the MSPE and store it
190 MSPE.ridge = get.MSPE(Y.valid, pred.ridge)
191 all.MSPEs[i, "Ridge"] = MSPE.ridge
192
193 matrix.train.raw = model.matrix(Ozone ~ ., data = data.train)
194 matrix.train = matrix.train.raw[,-1]
195
196 ### LASSO
197 all.LASSOs = cv.glmnet(x = matrix.train, y = Y.train)
198
199 ### Get both 'best' lambda values using $lambda.min and $lambda.1se
200 lambda.min = all.LASSOs$lambda.min
201 lambda.1se = all.LASSOs$lambda.1se

```

```

202
203 ### Get the coefficients for our two 'best' LASSO models
204 coef.LASSO.min = predict(all.LASSOs, s = lambda.min, type = "coef")
205 coef.LASSO.1se = predict(all.LASSOs, s = lambda.1se, type = "coef")
206
207 ### Get which predictors are included in our models (i.e. which
208 ### predictors have non-zero coefficients)
209 included.LASSO.min = predict(all.LASSOs, s = lambda.min,
210                             type = "nonzero")
211 included.LASSO.1se = predict(all.LASSOs, s = lambda.1se,
212                             type = "nonzero")
213
214 matrix.valid.LASSO.raw = model.matrix(Ozone ~ ., data = data.valid)
215 matrix.valid.LASSO = matrix.valid.LASSO.raw[,-1]
216 pred.LASSO.min = predict(all.LASSOs, newx = matrix.valid.LASSO,
217                          s = lambda.min, type = "response")
218 pred.LASSO.1se = predict(all.LASSOs, newx = matrix.valid.LASSO,
219                          s = lambda.1se, type = "response")
220
221 ### Calculate MSPEs and store them
222 MSPE.LASSO.min = get.MSPE(Y.valid, pred.LASSO.min)
223 all.MSPEs[i, "LASSO-Min"] = MSPE.LASSO.min
224
225 MSPE.LASSO.1se = get.MSPE(Y.valid, pred.LASSO.1se)
226 all.MSPEs[i, "LASSO-1se"] = MSPE.LASSO.1se
227
228 ## GAM
229 fit.gam = gam(Ozone ~ s(Solar.R) + s(Wind) + s(Temp) + s(TWcp) + s(TWrat),
230              data = data.train)
231
232 pred.gam = predict(fit.gam, data.valid)
233 MSPE.gam = get.MSPE(Y.valid, pred.gam) # Our helper function
234 all.MSPEs[i, "GAM"] = MSPE.gam
235
236
237 - #####
238 - ### PPR ###
239 - #####
240
241 ### To fit PPR, we need to do another round of CV. This time, do 5-fold
242 K.ppr = 5
243 n.train = nrow(data.train)
244 folds.ppr = get.folds(n.train, K.ppr)
245
246 ### Container to store MSPEs for each number of terms on each sub-fold
247 MSPEs.ppr = array(0, dim = c(max.terms, K.ppr))
248
249 - for(j in 1:K.ppr){
250   ### Split the training data.
251   ### Be careful! We are constructing an internal validation set by
252   ### splitting the training set from outer CV.
253   train.ppr = data.train[folds.ppr != j,]
254   valid.ppr = data.train[folds.ppr == j,]
255   Y.valid.ppr = valid.ppr$Ozone
256
257   ### We need to fit several different PPR models, one for each number
258   ### of terms. This means another for loop (make sure you use a different
259   ### index variable for each loop).
260   - for(l in 1:max.terms){
261     ### Fit model
262     fit.ppr = ppr(Ozone ~ ., data = train.ppr,
263                  max.terms = max.terms, nterms = l, sm.method = "gcv spline")
264
265     ### Get predictions and MSPE
266     pred.ppr = predict(fit.ppr, valid.ppr)
267     MSPE.ppr = get.MSPE(Y.valid.ppr, pred.ppr) # Our helper function
268
269     ### Store MSPE. Make sure the indices match for MSPEs.ppr
270     MSPEs.ppr[l, j] = MSPE.ppr
271   }
272 - }
273
274 ### Get average MSPE for each number of terms
275 ave.MSPE.ppr = apply(MSPEs.ppr, 1, mean)
276
277 ### Get optimal number of terms
278 best.terms = which.min(ave.MSPE.ppr)
279
280 ### Fit PPR on the whole CV training set using the optimal number of terms
281 fit.ppr.best = ppr(Ozone ~ ., data = data.train,
282                   max.terms = max.terms, nterms = best.terms, sm.method = "gcv spline")

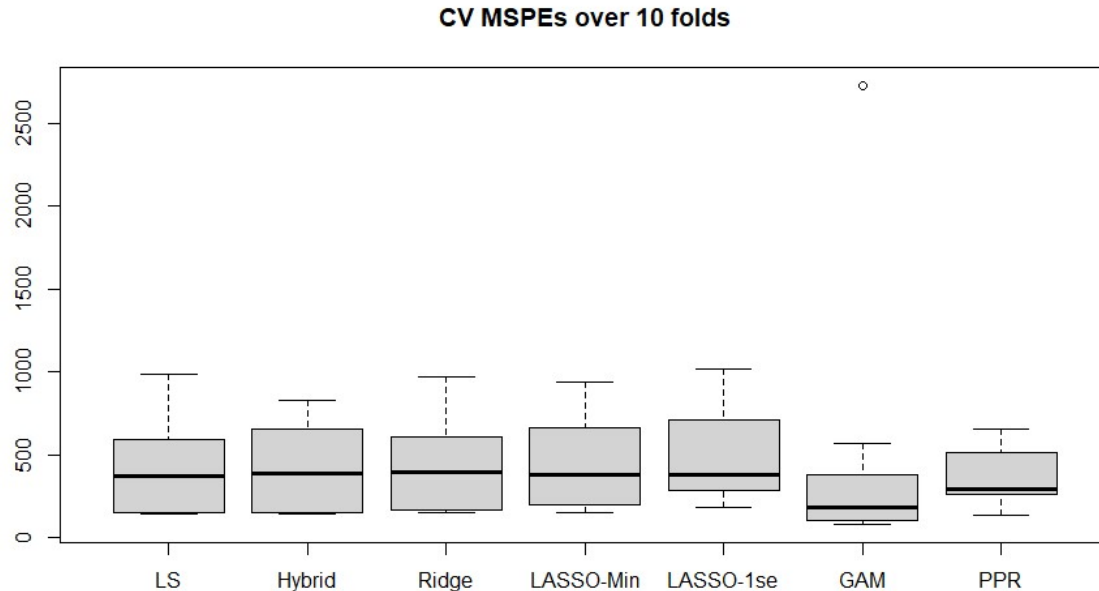
```

```

283
284   ### Get predictions, MSPE and store results
285   pred.ppr.best = predict(fit.ppr.best, data.valid)
286   MSPE.ppr.best = get.MSPE(Y.valid, pred.ppr.best) # Our helper function
287
288   all.MSPEs[i, "PPR"] = MSPE.ppr.best
289
290
291 }
292
293 all.MSPEs
294
295   ### Make a boxplot of MSPEs. I would like to include the number of folds
296   ### in the title. This can be done by using the paste0() function,
297   ### which concatenates strings (i.e. attaches them end-to-end), and
298   ### can be provided numeric variables.
299   boxplot(all.MSPEs, main = paste0("CV MSPEs over ", K, " folds"))
300
301
302
303   ### Calculate RMSPES
304   all.RMSPEs = apply(all.MSPEs, 1, function(w){
305     best = min(w)
306     return(w / best)
307   })
308   all.RMSPEs = t(all.RMSPEs)
309
310   ### Make a boxplot of RMSPES
311   boxplot(all.RMSPEs, main = paste0("CV RMSPES over ", K, " folds"))

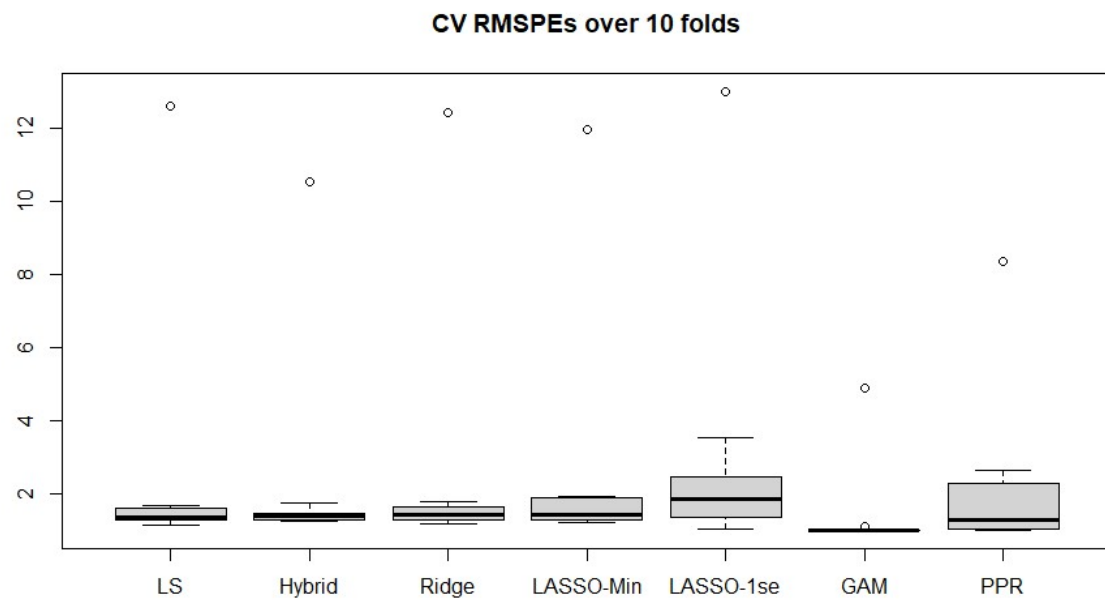
```

(a) Add the tuned PPR to the boxplots. Present the plots and write a sentence describing how well tuned PPR performs compared to other methods we have used thus far.



-> It looks like PPR is better than LASSO and Hybrid.

(b) Repeat this using relative MSPE.



(c) List the optimal tuning parameters that were selected for the tuned PPR in each of the 10 folds

```
[1] "10optimal number of terms5"
[1] "20optimal number of terms1"
[1] "30optimal number of terms1"
[1] "40optimal number of terms5"
[1] "50optimal number of terms1"
[1] "60optimal number of terms2"
[1] "70optimal number of terms1"
[1] "80optimal number of terms1"
[1] "90optimal number of terms1"
[1] "100optimal number of terms1"
```

5, 1, 1, 5, 1, 2, 1, 1, 1, 1