2. Add three new models fitted to all 5 variables to the 10-fold CV comparison that has been used for LASSO, Ridge, and other methods: Full tree (CP=0), min-cv-pruned tree, and 1-se-pruned tree. Use the same folds as before.

```r
120  get.folds = function(n, K) {
121    ### Get the appropriate number of fold labels
122    n.fold = ceiling(n / K) # Number of observations per fold (rounded up)
123    fold.ids.raw = rep(1:K, times = n.fold) # Generate extra labels
124    fold.ids = fold.ids.raw[1:n] # Keep only the correct number of labels
125
126    ### Shuffle the fold labels
127    folds.rand = fold.ids[sample.int(n)]
128    |
129    return(folds.rand)
130  }
131
132  # i=1
133  ### Number of folds
134  K = 10
135
136  ### Construct folds
137  n = nrow(data) # Sample size
138  folds = get.folds(n, K)
139
140  ### Create a container for MSPEs. Let's include ordinary least-squares
141  ### regression for reference
142  all.models = c("LS", "Hybrid", "Ridge", "LASSO-Min", "LASSO-1se", "GAM", "Full-tree", "Min-cv tree", "1-se tree")
143  all.MSPEs = array(0, dim = c(K, length(all.models)))
144  colnames(all.MSPEs) = all.models
145  ### Begin cross-validation
146  for(i in 1:K){
147    ### Split data
148    data.train = data[folds != i,]
149    data.valid = data[folds == i,]
150    n.train = nrow(data.train)
151
152    ### Get response vectors
153    Y.train = data.train$Ozone
154    Y.valid = data.valid$Ozone
155
156    # LS
157    fit.ls = lm(Ozone ~ ., data = data.train)
158    pred.ls = predict(fit.ls, newdata = data.valid)
159    MSPE.ls = get.MSPE(Y.valid, pred.ls)
160    all.MSPEs[i, "LS"] = MSPE.ls
161
162    #Hybrid Stepwise
```

```r
163
164    fit.start = lm(Ozone ~ 1, data = data.train)
165    fit.end = lm(Ozone ~ ., data = data.train)
166
167    step.BIC = step(fit.start, list(upper = fit.end), k = log(n.train),
168                    trace = 0)
169
170    pred.step.BIC = predict(step.BIC, data.valid)
171
172    err.step.BIC = get.MSPE(Y.valid, pred.step.BIC)
173
174    all.MSPEs[i, "Hybrid"] = err.step.BIC
175
176
177
178    #ridge regression
179    lambda.vals = seq(from = 0, to = 100, by = 0.05)
180
181    fit.ridge = lm.ridge(Ozone ~ ., lambda = lambda.vals,
182                         data = data.train)
183
184    ind.min.GCV = which.min(fit.ridge$GCV)
185    lambda.min = lambda.vals[ind.min.GCV]
186
187    all.coefs.ridge = coef(fit.ridge)
188    coef.min = all.coefs.ridge[ind.min.GCV,]
189
190    matrix.valid.ridge = model.matrix(Ozone ~ ., data = data.valid)
191
192    ### Now we can multiply the data by our coefficient vector. The
193    ### syntax in R for matrix-vector multiplication is %*%. Note that,
194    ### for this type of multiplication, order matters. That is,
195    ### A %*% B != B %*% A. Make sure you do data %*% coefficients.
196    ### For more information, see me in a Q&A session or, better still,
197    ### take a course on linear algebra (it's really neat stuff)
198    pred.ridge = matrix.valid.ridge %*% coef.min
199
200    ### Now we just need to calculate the MSPE and store it
201    MSPE.ridge = get.MSPE(Y.valid, pred.ridge)
202    all.MSPEs[i, "Ridge"] = MSPE.ridge
203
204    matrix.train.raw = model.matrix(Ozone ~ ., data = data.train)
205    matrix.train = matrix.train.raw[,-1]
206
```

```r
207    ### LASSO
208    all.LASSOs = cv.glmnet(x = matrix.train, y = Y.train)
209
210    ### Get both 'best' lambda values using $lambda.min and $lambda.1se
211    lambda.min = all.LASSOs$lambda.min
212    lambda.1se = all.LASSOs$lambda.1se
213
214    ### Get the coefficients for our two 'best' LASSO models
215    coef.LASSO.min = predict(all.LASSOs, s = lambda.min, type = "coef")
216    coef.LASSO.1se = predict(all.LASSOs, s = lambda.1se, type = "coef")
217
218    ### Get which predictors are included in our models (i.e. which
219    ### predictors have non-zero coefficients)
220    included.LASSO.min = predict(all.LASSOs, s = lambda.min,
221                                 type = "nonzero")
222    included.LASSO.1se = predict(all.LASSOs, s = lambda.1se,
223                                 type = "nonzero")
224
225    matrix.valid.LASSO.raw = model.matrix(Ozone ~ ., data = data.valid)
226    matrix.valid.LASSO = matrix.valid.LASSO.raw[,-1]
227    pred.LASSO.min = predict(all.LASSOs, newx = matrix.valid.LASSO,
228                             s = lambda.min, type = "response")
229    pred.LASSO.1se = predict(all.LASSOs, newx = matrix.valid.LASSO,
230                             s = lambda.1se, type = "response")
231
232    ### Calculate MSPEs and store them
233    MSPE.LASSO.min = get.MSPE(Y.valid, pred.LASSO.min)
234    all.MSPEs[i, "LASSO-Min"] = MSPE.LASSO.min
235
236    MSPE.LASSO.1se = get.MSPE(Y.valid, pred.LASSO.1se)
237    all.MSPEs[i, "LASSO-1se"] = MSPE.LASSO.1se
238
239    ## GAM
240    fit.gam = gam(Ozone ~ s(Solar.R) + s(Wind) + s(Temp) + s(TWcp) + s(TWrat),
241                  data = data.train)
242
243    pred.gam = predict(fit.gam, data.valid)
244    MSPE.gam = get.MSPE(Y.valid, pred.gam) # Our helper function
245    all.MSPEs[i, "GAM"] = MSPE.gam
246
247
248    # Full-tree
249    fit.tree = rpart(Ozone ~ ., data = data, cp=0)
```

```r
250    fit.tree.pred = predict(fit.tree, data.valid)
251    MSPE.fit.tree = get.MSPE(Y.valid, fit.tree.pred)
252    all.MSPEs[i, "Full-tree"] = MSPE.fit.tree
253
254
255    # Min-cv tree
256    info.tree = fit.tree$cptable
257    info.tree
258    ind.min = which.min(info.tree[,"xerror"])
259    CP.min.raw = info.tree[ind.min, "CP"]
260    if(ind.min == 1){
261      ### If minimum CP is in row 1, store this value
262      CP.min = CP.min.raw
263    } else{
264      ### If minimum CP is not in row 1, average this with the value from the
265      ### row above it.
266
267      ### Value from row above
268      CP.above = info.tree[ind.min-1, "CP"]
269
270      ### (Geometric) average
271      CP.min = sqrt(CP.min.raw * CP.above)
272    }
273    fit.tree.min = prune(fit.tree, cp = CP.min)
274    fit.tree.min.pred = predict(fit.tree.min, data.valid)
275    MSPE.fit.tree.min = get.MSPE(Y.valid, fit.tree.min.pred)
276    all.MSPEs[i, "Min-cv tree"] = MSPE.fit.tree.min
277
278
279    #"1-se tree"
280    err.min = info.tree[ind.min, "xerror"]
281    se.min = info.tree[ind.min, "xstd"]
282    threshold = err.min + se.min
283    ind.1se = min(which(info.tree[1:ind.min,"xerror"] < threshold))
284
285    ### Get the corresponding CP value, averaging if necessary
286    CP.1se.raw = info.tree[ind.1se, "xerror"]
287    if(ind.1se == 1){
288      ### If best CP is in row 1, store this value
289      CP.1se = CP.1se.raw
290    } else{
291      ### If best CP is not in row 1, average this with the value from the
292      ### row above it.
```

```
293
294        ### Value from row above
295        CP.above = info.tree[ind.1se-1, "CP"]
296
297        ### (Geometric) average
298        CP.1se = sqrt(CP.1se.raw * CP.above)
299 ~    }
300
301      fit.tree.1se = prune(fit.tree, cp = CP.1se)
302      fit.tree.1se.pred = predict(fit.tree.1se, data.valid)
303      MSPE.fit.tree.1se = get.MSPE(Y.valid, fit.tree.1se.pred)
304      all.MSPEs[i, "1-se tree"] = MSPE.fit.tree.1se
305 ~ }
306
307   all.MSPEs
308
309   ### Make a boxplot of MSPEs. I would like to include the number of folds
310   ### in the title. This can be done by using the paste0() function,
311   ### which concatenates strings (i.e. attaches them end-to-end), and
312   ### can be provided numeric variables.
313   boxplot(all.MSPEs, main = paste0("CV MSPEs over ", K, " folds"))
314
315
316
317   ### Calculate RMSPEs
318 ~ all.RMSPEs = apply(all.MSPEs, 1, function(W){
319      best = min(W)
320      return(W / best)
321 ~ })
322   all.RMSPEs = t(all.RMSPEs)
323
324   ### Make a boxplot of RMSPEs
325   boxplot(all.RMSPEs, main = paste0("CV RMSPEs over ", K, " folds"))
> all.MSPEs
          LS   Hybrid      Ridge LASSO-Min LASSO-1se        GAM Full-tree Min-cv tree 1-se tree
 [1,] 1037.2843 935.5934 1026.2063 1031.8530 1081.5330 135.46137 521.42750   531.63564 1103.9618
 [2,]  419.5374 520.7976  462.0808  463.0394  426.9839 1166.67426 129.56077   122.76918  598.1703
 [3,]  561.5317 606.2170  569.6294  559.5286  427.4034  674.21375 209.19630   209.19630  284.9489
 [4,]  117.0370 140.0399  119.3903  119.7927  194.3578  140.95548  98.93844    98.93844  199.0212
 [5,]  200.1188 195.1213  201.8795  230.7676  285.8414   42.07896 115.41266   112.12023  458.6532
 [6,]  208.1940 180.5834  208.5333  209.9717  220.0076  297.97684 386.37579   386.37579  752.4572
 [7,]  507.6231 555.5441  515.8408  654.4310 1062.9010  231.24568 158.20856   163.34328  305.6825
 [8,]  151.2711 150.9220  149.6126  150.8763  319.8533  199.88708 234.91171   254.13553  336.9521
 [9,]  681.4810 731.9910  687.8407  749.9537  928.2086  854.82298 418.22831   421.05592  434.6720
[10,]  358.1365 256.2422  365.3990  398.0162  281.7134  262.95811 218.79173   218.79173  242.9717
327   mean(all.MSPEs[,7])
328   mean(all.MSPEs[,8])
329   mean(all.MSPEs[,9])
```

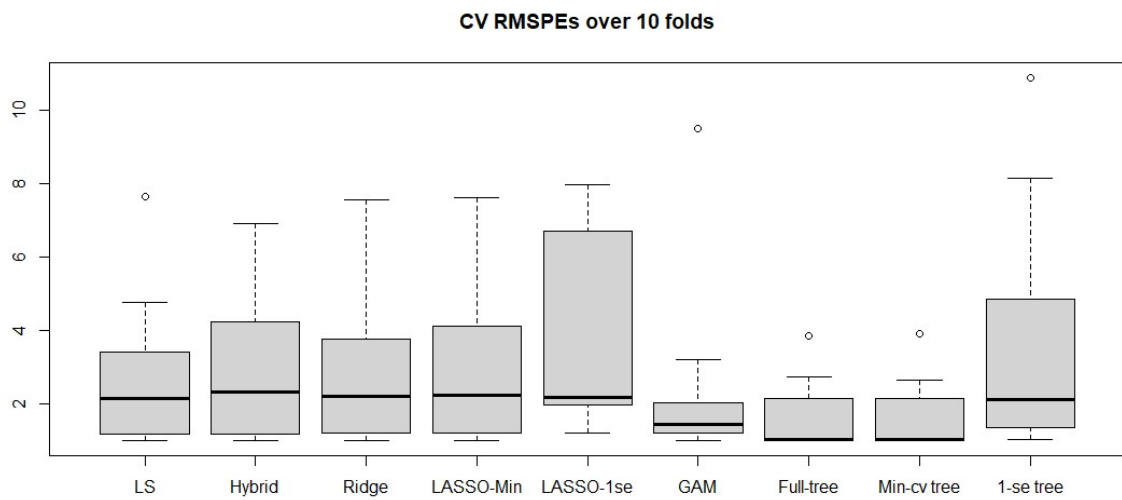(a) **Compute the mean MSPE for each tree model and comment on the comparison.**

```
> mean(all.MSPEs[,7])
[1] 249.1052
> mean(all.MSPEs[,8])
[1] 251.8362
> mean(all.MSPEs[,9])
[1] 471.7491
```

Full-tree and Min-cv tree has similar small MSPEs and 1-se tree has quite larger MSPE.

(b) **ADD the three trees to the relative MSPE boxplots made previously.**
**Comment on how well they perform compared to other methods.**

**CV RMSPEs over 10 folds**



So far, it looks like Full-tree and Min-cv tree and GAM show the best result. However, 1-se tree doesn't do well on it.