

OZONE DATA

Refer Lecture 6, Problem 3, where you used 10-fold CV for comparing LASSO models
Using *the same 10 folds* estimate the MSPE for PLS

```
get.folds = function(n, K) {  
  ### Get the appropriate number of fold labels  
  n.fold = ceiling(n / K) # Number of observations per fold (rounded up)  
  fold.ids.raw = rep(1:K, times = n.fold) # Generate extra labels  
  fold.ids = fold.ids.raw[1:n] # Keep only the correct number of labels  
  
  ### Shuffle the fold labels  
  folds.rand = fold.ids[sample.int(n)]  
  
  return(folds.rand)  
}  
  
### Number of folds  
K = 10  
  
### Construct folds  
data = ins.dv  
n = nrow(data) # Sample size  
folds = get.folds(n, K)  
  
### Create a container for MSPEs. Let's include ordinary least-squares  
### regression for reference  
all.models = c("PLS")  
all.MSPEs = array(0, dim = c(K, length(all.models)))  
colnames(all.MSPEs) = all.models
```

```

74 for(i in 1:K){
75   ### Split data
76   data.train = data[folds != i,]
77   data.valid = data[folds == i,]
78   n.train = nrow(data.train)
79
80   ### Get response vector
81   Y.valid = data.valid$per
82
83
84   fit.pls = plsr(per ~ ., data = data.train, validation = "CV",
85                 segments = 5)
86
87   #Code from L7 Question
88   mp.cv = fit.pls$validation
89   Opt.Comps = which.min(sqrt(mp.cv$PRESS/nrow(data)))
90   print(Opt.Comps)
91   ### Investigate the fitted PLS model. Comment out the next two
92   ### lines when running a CV loop
93
94   ### The summary function gives us lots of information about how
95   ### errors change as we increase the number of components
96   # summary(fit.pls)
97
98   ### The validationplot() function shows how MSPE from the internal
99   ### CV of plsr() changes with the number of included components.
100  # validationplot(fit.pls)
101
102  ### Get the best model from PLS. To do this, we need to find the model
103  ### that minimizes MSPE for the plsr() function's internal CV. It
104  ### takes a few steps, but all the information we need is contained
105  ### in the output of plsr().
106  CV.pls = fit.pls$validation # All the CV information
107  PRESS.pls = CV.pls$PRESS # Sum of squared CV residuals
108  CV.MSPE.pls = PRESS.pls / nrow(data.train) # MSPE for internal CV
109  ind.best.pls = which.min(CV.MSPE.pls) # Optimal number of components
110
111  get.MSPE = function(Y, Y.hat){
112    return(mean((Y - Y.hat)^2))
113  }
114
115  ### Get predictions and calculate MSPE on the validation fold
116  ### Set ncomps equal to the optimal number of components
117  pred.pls = predict(fit.pls, data.valid, ncomp = ind.best.pls)
118  MSPE.pls = get.MSPE(Y.valid, pred.pls)
119  all.MSPEs[i, "PLS"] = MSPE.pls
120 }

```

(a) On each training set, run the PLS function, choosing the optimum number of components using 10-fold CV. You can figure out how many components is optimal according to CV using code like this, assuming that I have named my model,

```
mod.pls
```

```
mp.cv = mod.pls$validation
```

```
Opt.Comps = which.min(sqrt(mp.cv$PRESS/nrow(ins.dv)))
```

Report the optimal number of components for each of the 10 folds.

```
[1] 10
[1] 11
[1] 7
[1] 8
[1] 8
[1] 8
[1] 8
[1] 7
[1] 6
[1] 7
```

(b) Use the predict() function to predict responses on the test set. You need to specify ncomps for this function.

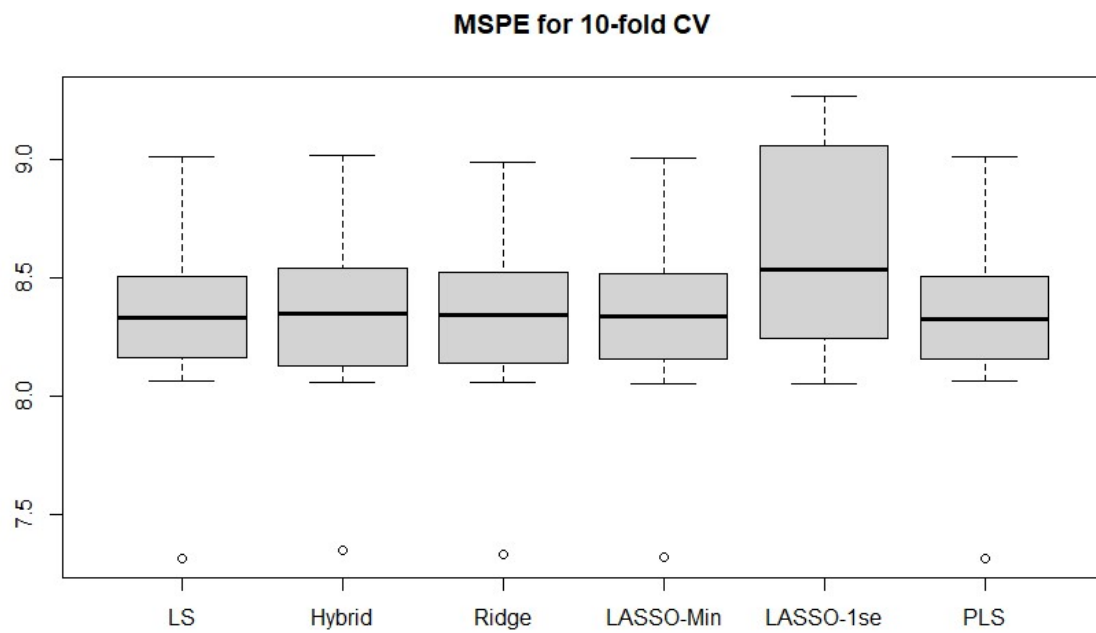
(c) **Report the separate MSPEs from each fold, $MSPE_v$, $v = 1, \dots, 10$ and the MSPE for the full data.**

```
> all.MSPEs
```

```
          PLS
[1,] 8.189934
[2,] 7.309847
[3,] 8.505827
[4,] 8.515403
[5,] 8.190235
[6,] 8.162715
[7,] 8.541857
[8,] 8.063142
[9,] 8.509003
[10,] 9.009496
```

(d) **ADD a boxplot for PLS to the boxplots you made for other models.**

Comment on how well PLS compares to the other models.



-> PLS does pretty good because it seems to have the lowest MSPE value.

(e) Remake the plot using relative MSPE and comment.

