

## 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).

Use regression trees to model the relationship between Ozone and all five explanatories as specified below.:

1. Fit a regression tree using only Temp and Wind, using  $cp=0$ .

(a) Print out the  $cp$  table.

```
1 # Title: STAT 452 Exercise 8 L13Q1
2 # Author: Injun Son
3 # Date: October 31, 2020
4
5 library(dplyr)
6 library(MASS) # For ridge regression
7 library(glmnet) # For LASSO
8 library(nnet) # Fits neural net models
9 library(rgl)
10 library(rpart)
11 library(rpart.plot)
12 source("Helper Functions.R")
13 data = na.omit(airquality[, 1:4])
14 data$Twcp = data$Temp*data$Wind
15 data$Twrat = data$Temp/data$Wind
16
17 set.seed(8598176)
18
19 # 1. Fit a regression tree using only Temp and Wind, using  $cp=0$ .
20 # (a) Print out the  $cp$  table.
21 # i. According to minimum CV, what value of  $cp$  should be used for optimal
22 # pruning, and what is the relative CV error at this value?
23
24 fit.tree = rpart(Ozone ~ Temp + Wind, data = data, cp=0)
25
26 ### Get the CP table
27 info.tree = fit.tree$cptable
```

	CP	nsplit	rel error	xerror	xstd
1	0.484385708	0	1.0000000	1.0068193	0.1705465
2	0.097982838	1	0.5156143	0.5469021	0.1798443
3	0.057062146	2	0.4176315	0.5690217	0.1811551
4	0.019726902	3	0.3605693	0.5260658	0.1525831
5	0.018716394	4	0.3408424	0.5185840	0.1519412
6	0.005142355	5	0.3221260	0.4861408	0.1409565
7	0.004575931	6	0.3169837	0.4878673	0.1409443
8	0.001882628	7	0.3124077	0.4930603	0.1410539
9	0.000000000	9	0.3086425	0.4901064	0.1410692

i. According to minimum CV, **what value of cp should be used for optimal pruning, and what is the relative CV error at this value?**

```
30 ### We have to prune the tree manually. First, get the CP value with minimum
31 ### CV error
32 ind.min = which.min(info.tree[, "xerror"])
33 CP.min.raw = info.tree[ind.min, "CP"]
```

```
> CP.min.raw
```

```
[1] 0.005142355
```

Relative CV error is 0.3221260

(b) **Plot both the min-CV pruned tree and the 1SE pruned tree.**

i. Interpret the first split on the root node: **What variable and split location are used, how many observations go each way, and what are the new means in each node?**

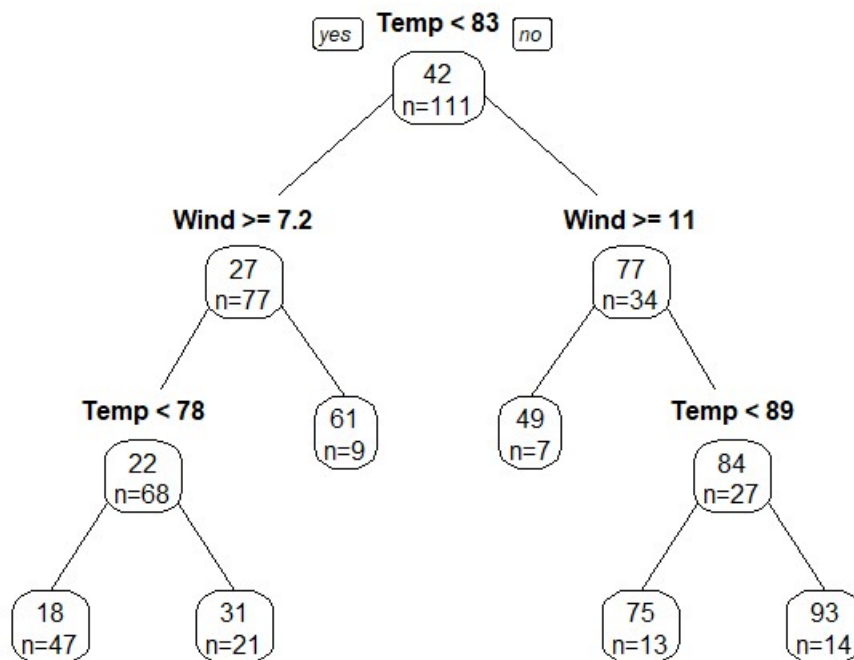
```
36 #####
37 # (b) Plot both the min-CV pruned tree and the 1SE pruned tree.
38 # i. Interpret the first split on the root node: What variable and split location
39 # are used, how many observations go each way, and what are
40 # the new means in each node?
41
42 ### It's best to average the minimum CP value with the one from the row above
43 ### using the geometric mean (i.e. multiply them together, then square root).
44 ### If we implement this procedure directly, and the minimum CP value is in
45 ### the first row, our code will probably give an error. We should write our
46 ### code so that it can cope with this weird situation. This attitude is
47 ### called defensive programming, and it is a very good habit to practice.
48
49 ### Check if minimum CP value is in row 1. We can do this using an if-else
50 ### statement. See this tutorial's video for details.
51 if(ind.min == 1){
52   ### If minimum CP is in row 1, store this value
53   CP.min = CP.min.raw
54 } else{
55   ### If minimum CP is not in row 1, average this with the value from the
56   ### row above it.
57
58   ### Value from row above
59   CP.above = info.tree[ind.min-1, "CP"]
60
61   ### (Geometric) average
62   CP.min = sqrt(CP.min.raw * CP.above)
63 }
64
65 ### We now have our chosen CP value. We can prune our tree using the prune()
66 ### function. The first input to prune() is the tree object, and we set "cp"
67 ### equal to the CP value where we want to prune
68 fit.tree.min = prune(fit.tree, cp = CP.min)
69
70 ### Next, we want to prune using the 1SE rule. Fortunately, the CP table
71 ### gives us the CV standard error. First, find the minimum CV error plus 1
72 ### standard error
73 err.min = info.tree[ind.min, "xerror"]
74 se.min = info.tree[ind.min, "xstd"]
75 threshold = err.min + se.min
76
```

```

77  ### Next, get the smallest tree with CV error below our threshold.
78  ### Note: We limit our search to only trees which are no larger than our min CV
79  ###      error tree.
80  ind.lse = min(which(info.tree[1:ind.min,"xerror"] < threshold))
81
82  ### Get the corresponding CP value, averaging if necessary
83  CP.lse.raw = info.tree[ind.lse, "xerror"]
84  if(ind.lse == 1){
85      ### If best CP is in row 1, store this value
86      CP.lse = CP.lse.raw
87  } else{
88      ### If best CP is not in row 1, average this with the value from the
89      ### row above it.
90
91      ### Value from row above
92      CP.above = info.tree[ind.lse-1, "CP"]
93
94      ### (Geometric) average
95      CP.lse = sqrt(CP.lse.raw * CP.above)
96  }
97
98  ### Prune the tree
99  fit.tree.lse = prune(fit.tree, cp = CP.lse)
100
101  #####
102  ### A nice feature of tree models is that they make great plots. Let's ###
103  ### plot the full tree, and both pruned trees. We can plot trees using ###
104  ### the prp() function from the rpart.plot package. ###
105  #####
106  ### The prp() function has many inputs We will just use two: type and extra.
107  ### Setting both of these inputs to 1 gives a nice looking plot. Since
108  ### prp() makes a plot, we can set the title using main. We also
109  ### have to provide the fitted tree object which is being plotted.
110  prp(fit.tree, type = 1, extra = 1, main = "Full Tree")
111  prp(fit.tree.min, type = 1, extra = 1, main = "Pruned Tree - Min")
112  prp(fit.tree.lse, type = 1, extra = 1, main = "Pruned Tree - lse")

```

## Pruned Tree - Min



```
> fit.tree.min
```

```
n= 111
```

```
node), split, n, deviance, yval
* denotes terminal node
```

- 1) root 111 121801.900 42.09910
- 2) Temp< 82.5 77 42143.250 26.77922
- 4) wind>=7.15 68 10886.750 22.25000
- 8) Temp< 77.5 47 3863.404 18.27660 \*
- 9) Temp>=77.5 21 4620.571 31.14286 \*
- 5) wind< 7.15 9 19322.000 61.00000 \*
- 3) Temp>=82.5 34 20659.560 76.79412
- 6) wind>=10.6 7 1183.429 48.71429 \*
- 7) wind< 10.6 27 12525.850 84.07407
- 14) Temp< 88.5 13 7307.231 74.53846 \*
- 15) Temp>=88.5 14 2938.929 92.92857 \*

Temp variable was used in the root node and 77 variables that Temp<82.5 goes left and 34 variables that Temp >=82.5 goes right. The new mean is 26.77 for variables that Temp< 82.5 and the new mean is 76.79

### Pruned Tree - 1SE

42  
n=111

```
> fit.tree.1se  
n= 111  
  
node), split, n, deviance, yval  
    * denotes terminal node  
  
1) root 111 121801.9 42.0991 *
```

1SE pruned tree didn't split the data.