

State Data Revisted

Contents

Introduction	2
Exercices	2
1. Linear Regrerssion Models	2
Problem 1.1	2
Problem 1.2	3
Problem 1.3	3
Problem 1.4	4
Problem 1.5	4
2. CART Models	5
Problem 2.1	5
Problem 2.2	6
Problem 2.3	6
Problem 2.4	7
Problem 2.5	8
Problem 2.6	8
Problem 2.7	8
3. Cross-Validation	9
Problem 3.1	9
Problem 3.2	10
Problem 3.3	11
Problem 3.4	12
Problem 3.5	12
Problem 3.6	14
Problem 3.7	15

Introduction

We will be revisiting the “state” dataset from one of the optional problems in Unit 2. This dataset has, for each of the fifty U.S. states, the population, per capita income, illiteracy rate, murder rate, high school graduation rate, average number of frost days, area, latitude and longitude, division the state belongs to, region the state belongs to, and two-letter abbreviation. This dataset comes from the U.S. Department of Commerce, Bureau of the Census.

Load the dataset into R and convert it to a data frame by running the following two commands in R:

If you can’t access the state dataset in R, here is a CSV file with the same data that you can load into R using the `read.csv` function: `statedataSimple.csv`. Be sure to call the output of the `read.csv` function “`statedata`”.

After you have loaded the data into R, inspect the data set using the command: `str(statedata)`

This dataset has 50 observations (one for each US state) and the following 8 variables:

Population : the population estimate of the state in 1975

Income : per capita income in 1974

Illiteracy : illiteracy rates in 1970, as a percent of the population

Life.Exp : the life expectancy in years of residents of the state in 1970

Murder : the murder and non-negligent manslaughter rate per 100,000 population in 1976

HS.Grad : percent of high-school graduates in 1970

Frost : the mean number of days with minimum temperature below freezing from 1931 - 1960 in the capital or a large city of the state

Area : the land area (in square miles) of the state

We will try to build a model for life expectancy using regression trees, and employ cross-validation to improve our tree’s performance.

Exercices

1. Linear Regrerssion Models

Problem 1.1 Let’s recreate the **linear regression** models we made in the previous homework question. First, predict `Life.Exp` using all of the other variables as the independent variables (`Population`, `Income`, `Illiteracy`, `Murder`, `HS.Grad`, `Frost`, `Area`). Use the entire dataset to build the model.

What is the adjusted R-squared of the model?

```
## 'data.frame':   50 obs. of  8 variables:
## $ Population: int  3615 365 2212 2110 21198 2541 3100 579 8277 4931 ...
## $ Income    : int  3624 6315 4530 3378 5114 4884 5348 4809 4815 4091 ...
## $ Illiteracy: num  2.1 1.5 1.8 1.9 1.1 0.7 1.1 0.9 1.3 2 ...
## $ Life.Exp  : num  69 69.3 70.5 70.7 71.7 ...
## $ Murder    : num  15.1 11.3 7.8 10.1 10.3 6.8 3.1 6.2 10.7 13.9 ...
## $ HS.Grad   : num  41.3 66.7 58.1 39.9 62.6 63.9 56 54.6 52.6 40.6 ...
## $ Frost     : int   20 152 15 65 20 166 139 103 11 60 ...
## $ Area      : int  50708 566432 113417 51945 156361 103766 4862 1982 54090 58073 ...

##
## Call:
## lm(formula = Life.Exp ~ ., data = statedata)
##
## Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -1.48895 -0.51232 -0.02747  0.57002  1.49447
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.094e+01  1.748e+00  40.586 < 2e-16 ***
## Population   5.180e-05  2.919e-05   1.775  0.0832 .
## Income      -2.180e-05  2.444e-04  -0.089  0.9293
## Illiteracy   3.382e-02  3.663e-01   0.092  0.9269
## Murder      -3.011e-01  4.662e-02  -6.459 8.68e-08 ***
## HS.Grad      4.893e-02  2.332e-02   2.098  0.0420 *
## Frost       -5.735e-03  3.143e-03  -1.825  0.0752 .
## Area        -7.383e-08  1.668e-06  -0.044  0.9649
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7448 on 42 degrees of freedom
## Multiple R-squared:  0.7362, Adjusted R-squared:  0.6922
## F-statistic: 16.74 on 7 and 42 DF,  p-value: 2.534e-10
```

Answer : 0.6922

Explanation :

To build the regression model, type the following command in your R console:

Then, if you look at the output of `summary(RegModel)`, you should see that the Adjusted R-squared is 0.6922.

Problem 1.2 Calculate the sum of squared errors (SSE) between the predicted life expectancies using this model and the actual life expectancies:

$$SSE = \sum (PredictValue - TestValue)^2$$

```
## [1] 23.29714
```

Answer : 23.29714

Explanation :

To make predictions, type in your R console:

where “RegModel” is the name of your regression model. You can then compute the sum of squared errors by typing the following in your R console:

The SSE is 23.29714.

Alternatively, you can use the following command to get the SSE:

Problem 1.3 Build a second linear regression model using just **Population**, **Murder**, **Frost**, and **HS.Grad** as independent variables (the best 4 variable model from the previous homework).

What is the adjusted R-squared for this model?

```
##
## Call:
## lm(formula = Life.Exp ~ Population + Murder + Frost + HS.Grad,
##     data = statedata)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.47095 -0.53464 -0.03701  0.57621  1.50683
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.103e+01  9.529e-01  74.542 < 2e-16 ***
## Population    5.014e-05  2.512e-05   1.996  0.05201 .
## Murder       -3.001e-01  3.661e-02  -8.199 1.77e-10 ***
## Frost        -5.943e-03  2.421e-03  -2.455  0.01802 *
## HS.Grad       4.658e-02  1.483e-02   3.142  0.00297 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7197 on 45 degrees of freedom
## Multiple R-squared:  0.736, Adjusted R-squared:  0.7126
## F-statistic: 31.37 on 4 and 45 DF,  p-value: 1.696e-12
```

Answer : 0.7126

Explanation :

You can create this regression model by typing the following into your R console:

The Adjusted R-squared is at the bottom right of the output, and is 0.7126

Problem 1.4 Calculate the sum of squared errors again, using this reduced model:

Answer : 23.30804

Explanation :

The sum of squared errors (SSE) can be computed by first making predictions:

and then computing the sum of the squared difference between the actual values and the predictions:

Alternatively, you can compute the SSE with the following command:

Problem 1.5 Which of the following is correct?

Answer :

1. **Trying different combinations of variables in linear regression is like trying different numbers of splits in a tree - this controls the complexity of the model.**
2. Using many variables in a linear regression is always better than using just a few.
3. The variables we removed were uncorrelated with Life.Exp

Explanation :

The correct answer is the first one. Trying different combinations of variables in linear regression controls the complexity of the model. This is similar to trying different numbers of splits in a tree, which is also controlling the complexity of the model.

The second answer is incorrect because as we see here, a model with fewer variables actually has a higher adjusted R-squared. If your accuracy is just as good, a model with fewer variables is almost always better.

The third answer is incorrect because the variables we removed have non-zero correlations with the dependent variable Life.Exp. Illiteracy and Area are negatively correlated with Life.Exp, with correlations of -0.59 and -0.11. Income is positively correlated with Life.Exp, with a correlation of 0.34. These correlations can be computed by typing the following into your R console:

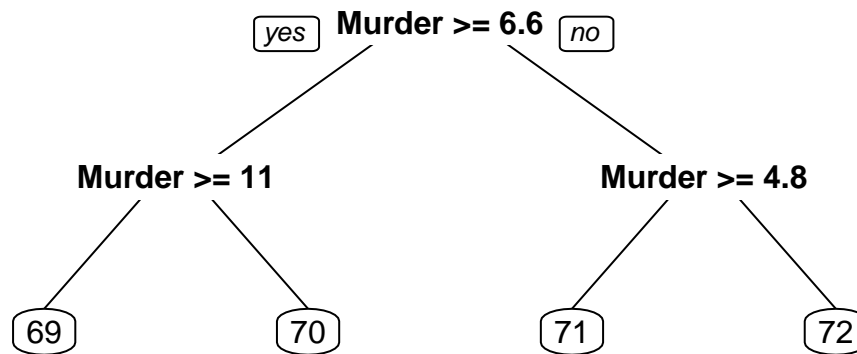
2. CART Models

Problem 2.1 Let's now build a CART model to predict Life.Exp using all of the other variables as independent variables (Population, Income, Illiteracy, Murder, HS.Grad, Frost, Area). We'll use the default minbucket parameter, so don't add the minbucket argument. Remember that in this problem we are not as interested in predicting life expectancies for new observations as we are understanding how they relate to the other variables we have, so we'll use all of the data to build our model. You shouldn't use the method="class" argument since this is a regression tree.

Plot the tree.

Which of these variables appear in the tree?

Select all that apply.



Answer :

1. Population
2. **Murder**
3. Frost

4. HS.Grad
5. Area

Explanation :

You can create the tree in R by typing the following command:

Be sure to load the “rpart” and “rpart.plot” packages with the library command if they are not already loaded.

You can then plot the tree by typing:

We can see that the only variable used in the tree is “Murder”.

Problem 2.2 Use the regression tree you just built to predict life expectancies (using the predict function), and calculate the sum-of-squared-errors (SSE) like you did for linear regression.

What is the SSE?

$$SSE = \sum (PredictValue - TestValue)^2$$

```
## [1] 28.99848
```

Answer : 28.99848

Explanation :

You can make predictions using the CART model by typing the following line in your R console (assuming your model is called “CARTmodel”):

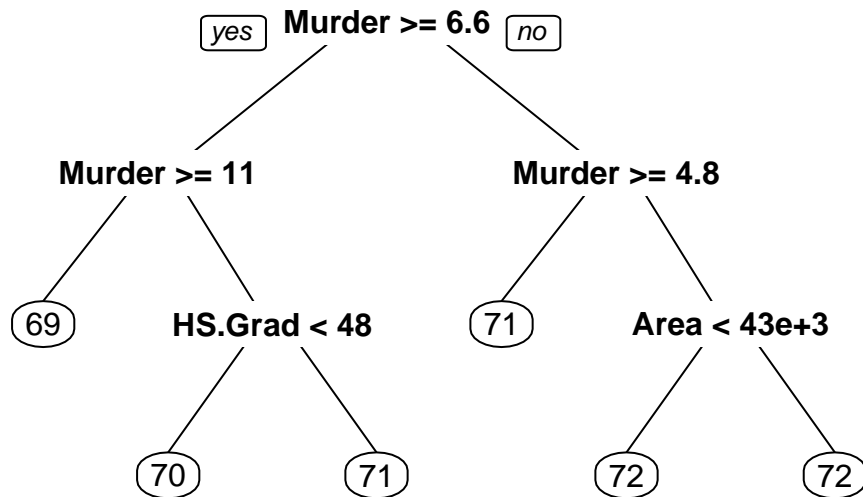
Then, you can compute the sum of squared errors (SSE) by typing the following:

The SSE is 28.99848.

Problem 2.3 The error is higher than for the linear regression models. One reason might be that we haven’t made the tree big enough. Set the minbucket parameter to 5, and recreate the tree.

Which variables appear in this new tree?

Select all that apply.



Answer :

1. Population
2. **Murder** correct
3. Frost
4. **HS.Grad**
5. **Area**

Explanation :

You can create a tree with a minbucket value of 5 with the following command:

Then, if you plot the tree using

you can see that Murder, HS.Grad, and Area are all used in this new tree.

Problem 2.4 Do you think the default minbucket parameter is smaller or larger than 5 based on the tree that was built?

Answer :

1. Smaller
2. **Larger**

Explanation :

Since the tree now has more splits, it must be true that the default minbucket parameter was limiting the tree from splitting more before. So the default minbucket parameter must be larger than 5.

Problem 2.5 What is the SSE of this tree?

```
## [1] 23.64283
```

Answer : 23.64283

Explanation :

You can compute the SSE of this tree by first making predictions:

and then computing the sum of the squared differences between the actual values and the predicted values:

The SSE is 23.64283

This is much closer to the linear regression model's error. By changing the parameters we have improved the fit of our model.

Problem 2.6 Can we do even better? Create a tree that predicts Life.Exp using **only Area**, with the minbucket parameter to 1.

What is the SSE of this newest tree?

```
## [1] 9.312442
```

Answer : 9.312442

Explanation :

You can create this third tree by typing:

Then to compute the SSE, first make predictions:

And then compute the sum of squared differences between the actual values and the predicted values:

The SSE is 9.312442.

Note that the SSE is not zero here - we still make some mistakes. This is because there are other parameters in rpart that are also trying to prevent the tree from overfitting by setting default values. So our tree doesn't necessarily have one observation in each bucket - by setting minbucket=1 we are just allowing the tree to have one observation in each bucket.

Problem 2.7 This is the lowest error we have seen so far.

What would be the best interpretation of this result?

Answer :

1. Trees are much better than linear regression for this problem because they can capture nonlinearities that linear regression misses.
2. **We can build almost perfect models given the right parameters, even if they violate our intuition of what a good model should be.**
3. Area is obviously a very meaningful predictor of life expectancy, given we were able to get such low error using just Area as our independent variable.

Explanation :

The correct answer is the second one. By making the minbucket parameter very small, we could build an almost perfect model using just one variable, that is not even our most significant variable. However, if you plot the tree using `prp(CARTmodel3)`, you can see that the tree has 22 splits! This is not a very interpretable model, and will not generalize well.

The first answer is incorrect because our tree model that was not overfit performed similarly to the linear regression model. Trees only look better than linear regression here because we are overfitting the model to the data.

The third answer is incorrect because Area is not actually a very meaningful predictor. Without overfitting the tree, our model would not be very accurate only using Area.

3. Cross-Validation

Problem 3.1 Adjusting the variables included in a linear regression model is a form of model tuning. In Problem 1 we showed that by removing variables in our linear regression model (tuning the model), we were able to maintain the fit of the model while using a simpler model. A rule of thumb is that simpler models are more interpretable and generalizable. We will now tune our regression tree to see if we can improve the fit of our tree while keeping it as simple as possible.

Load the caret library, and **set the seed to 111**. Set up the controls exactly like we did in the lecture (**10-fold cross-validation**) with **cp varying over the range 0.01 to 0.50 in increments of 0.01**. Use the train function to determine the best cp value for a CART model using all of the available independent variables, and the entire dataset statedata.

What value of cp does the train function recommend?

(Remember that the train function tells you to pick the largest value of cp with the lowest error when there are ties, and explains this at the bottom of the output.)

```
## CART
##
## 50 samples
## 7 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 44, 45, 45, 46, 44, 45, ...
## Resampling results across tuning parameters:
##
##   cp      RMSE      Rsquared    MAE
##   0.01  1.042909  0.5206939  0.8299237
##   0.02  1.042909  0.5206939  0.8299237
##   0.03  1.027748  0.5338091  0.8173654
##   0.04  1.032567  0.5338091  0.8173654
##   0.05  1.032567  0.5338091  0.8173654
##   0.06  1.033866  0.5206565  0.8320776
##   0.07  1.029830  0.5285854  0.8252412
##   0.08  1.029830  0.5285854  0.8252412
##   0.09  1.029830  0.5285854  0.8252412
##   0.10  1.005814  0.5512315  0.8079269
##   0.11  1.005814  0.5512315  0.8079269
##   0.12  1.005814  0.5512315  0.8079269
##   0.13  1.032234  0.5238042  0.8262453
##   0.14  1.083214  0.5041955  0.8725504
##   0.15  1.106834  0.4822947  0.9050228
##   0.16  1.138118  0.4775423  0.9422217
##   0.17  1.174001  0.4287787  0.9676503
##   0.18  1.192122  0.3990629  0.9942598
##   0.19  1.192122  0.3990629  0.9942598
##   0.20  1.192122  0.3990629  0.9942598
```

```
## 0.21 1.192122 0.3990629 0.9942598
## 0.22 1.192122 0.3990629 0.9942598
## 0.23 1.192122 0.3990629 0.9942598
## 0.24 1.192122 0.3990629 0.9942598
## 0.25 1.192122 0.3990629 0.9942598
## 0.26 1.192122 0.3990629 0.9942598
## 0.27 1.192122 0.3990629 0.9942598
## 0.28 1.192122 0.3990629 0.9942598
## 0.29 1.192122 0.3990629 0.9942598
## 0.30 1.192122 0.3990629 0.9942598
## 0.31 1.192122 0.3990629 0.9942598
## 0.32 1.192122 0.3990629 0.9942598
## 0.33 1.192122 0.3990629 0.9942598
## 0.34 1.192122 0.3990629 0.9942598
## 0.35 1.192122 0.3990629 0.9942598
## 0.36 1.192122 0.3990629 0.9942598
## 0.37 1.192122 0.3990629 0.9942598
## 0.38 1.192122 0.3990629 0.9942598
## 0.39 1.192122 0.3990629 0.9942598
## 0.40 1.192122 0.3990629 0.9942598
## 0.41 1.192122 0.3990629 0.9942598
## 0.42 1.192122 0.3990629 0.9942598
## 0.43 1.192122 0.3990629 0.9942598
## 0.44 1.192122 0.3990629 0.9942598
## 0.45 1.192122 0.3990629 0.9942598
## 0.46 1.308759 0.3091923 1.0963695
## 0.47 1.309534 0.3328614 1.0922317
## 0.48 1.358580 0.2908310 1.1126138
## 0.49 1.335777 0.3552299 1.0771938
## 0.50 1.361946 0.2358921 1.1229974
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was cp = 0.12.
```

Answer : 0.12

Explanation :

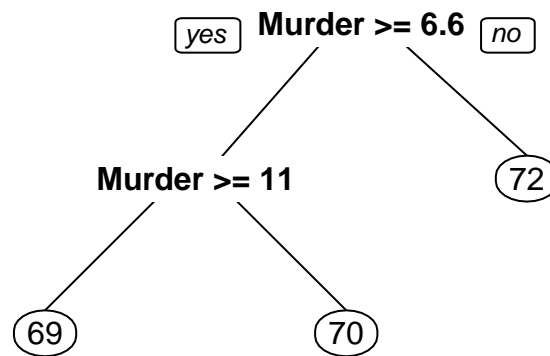
You can load the library caret and set the seed by typing the following commands:

Then, you can set up the cross-validation controls by typing:

You can then use the train function to find the best value of cp by typing:

At the bottom of the output, it says that the best value of cp is 0.12.

Problem 3.2 Create a tree with the value of cp you found in the previous problem, all of the available independent variables, and the entire dataset “statedata” as the training data. Then plot the tree.



You'll notice that this is actually quite similar to the first tree we created with the initial model. Interpret the tree: we predict the life expectancy to be 70 if the murder rate is greater than or equal to :

Answer : 6.6

and is less than :

Answer : 11

Explanation :

You can create a new tree with $cp=0.12$ by typing:

Then, if you plot the tree using

you can see that the life expectancy is predicted to be 70 if Murder is greater than or equal to 6.6 (the first split) and less than 11 (the second split).

Problem 3.3 Calculate the SSE of this tree:

[1] 32.86549

Answer : 32.86549

Explanation :

To compute the SSE, first make predictions:

and then compute the sum of squared differences between the actual values and the predicted values:

The SSE is 32.86549

Problem 3.4 Recall the first tree (default parameters), second tree (minbucket = 5), and the third tree (selected with cross validation) we made. Given what you have learned about cross-validation, **which of the three models would you expect to be better if we did use it for prediction on a test set?** For this question, suppose we had actually set aside a few observations (states) in a test set, and we want to make predictions on those states.

Answer :

1. The first model
2. The second model
3. **The model we just made with the “best” cp**

Explanation :

The purpose of cross-validation is to pick the tree that will perform the best on a test set. So we would expect the model we made with the “best” cp to perform best on a test set.

Problem 3.5 At the end of Problem 2 we made a very complex tree using just Area. Use train with the same parameters as before but just using Area as an independent variable to find the best cp value (set the seed to 111 first). Then build a new tree using just Area and this value of cp.

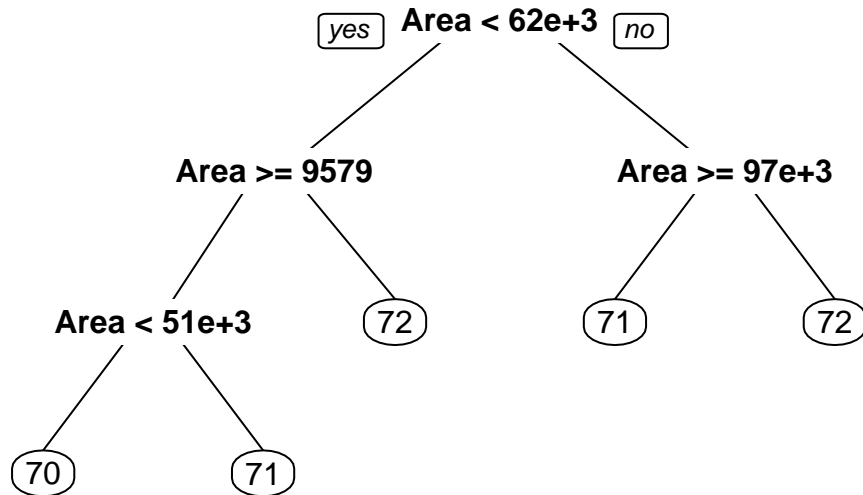
How many splits does the tree have?

```
## CART
##
## 50 samples
## 1 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 44, 45, 45, 46, 44, 45, ...
## Resampling results across tuning parameters:
##
##   cp      RMSE      Rsquared    MAE
##   0.01  1.285085  0.24288993  1.068979
##   0.02  1.297420  0.227958081  1.077964
##   0.03  1.297420  0.227958081  1.077964
##   0.04  1.297420  0.227958081  1.077964
##   0.05  1.297420  0.227958081  1.077964
##   0.06  1.283241  0.255224868  1.071996
##   0.07  1.283241  0.255224868  1.071996
##   0.08  1.277535  0.253025061  1.054684
##   0.09  1.286127  0.239816630  1.060619
##   0.10  1.286127  0.239816630  1.060619
##   0.11  1.286127  0.239816630  1.060619
##   0.12  1.278550  0.239816630  1.060619
##   0.13  1.336117  0.205007172  1.116064
##   0.14  1.364618  0.132092640  1.125916
##   0.15  1.364016  0.272311296  1.124202
##   0.16  1.348422  0.216040174  1.128918
##   0.17  1.365452  0.125811897  1.110581
##   0.18  1.334937  0.006222148  1.103295
##   0.19  1.328891           NaN  1.099663
##   0.20  1.328891           NaN  1.099663
##   0.21  1.328891           NaN  1.099663
```

```

## 0.22 1.328891      NaN 1.099663
## 0.23 1.328891      NaN 1.099663
## 0.24 1.328891      NaN 1.099663
## 0.25 1.328891      NaN 1.099663
## 0.26 1.328891      NaN 1.099663
## 0.27 1.328891      NaN 1.099663
## 0.28 1.328891      NaN 1.099663
## 0.29 1.328891      NaN 1.099663
## 0.30 1.328891      NaN 1.099663
## 0.31 1.328891      NaN 1.099663
## 0.32 1.328891      NaN 1.099663
## 0.33 1.328891      NaN 1.099663
## 0.34 1.328891      NaN 1.099663
## 0.35 1.328891      NaN 1.099663
## 0.36 1.328891      NaN 1.099663
## 0.37 1.328891      NaN 1.099663
## 0.38 1.328891      NaN 1.099663
## 0.39 1.328891      NaN 1.099663
## 0.40 1.328891      NaN 1.099663
## 0.41 1.328891      NaN 1.099663
## 0.42 1.328891      NaN 1.099663
## 0.43 1.328891      NaN 1.099663
## 0.44 1.328891      NaN 1.099663
## 0.45 1.328891      NaN 1.099663
## 0.46 1.328891      NaN 1.099663
## 0.47 1.328891      NaN 1.099663
## 0.48 1.328891      NaN 1.099663
## 0.49 1.328891      NaN 1.099663
## 0.50 1.328891      NaN 1.099663
##
## RMSE was used to select the optimal model using the smallest value.
## The final value used for the model was cp = 0.08.

```



Answer : 4

Explanation :

To find the best value of `cp` when using only `Area`, use the following command:

Then, build a new CART tree by typing:

You can plot the tree with `prp(CARTmodel5)`, and see that the tree has 4 splits.

Problem 3.6 The lower left leaf (or bucket) corresponds to the lowest predicted `Life.Exp` of 70. Observations in this leaf correspond to states with area greater than or equal to :

Answer : 9579

and area less than :

Answer : 51000

Explanation :

To get to this leaf, we go through 3 splits:

Area less than 62,000

Area greater than or equal to 9,579

Area less than 51,000

This means that this leaf is composed of states that have an area greater than 9,579 and less than 51,000.

Problem 3.7 We have simplified the previous “Area tree” considerably by using cross-validation. **Calculate the SSE of the cross-validated “Area tree”, and select all of the following correct statements that apply:**

[1] 44.26817

Answer :

1. The best model in this whole question is the first “Area tree” because it had the lowest SSE.
2. **The Area variable is not as predictive as Murder rate.**
3. Cross-validation is intended to decrease the SSE for a model on the training data, compared to a tree that isn’t cross-validated.
4. Cross-validation will always improve the SSE of a model on unseen data, compared to a tree that isn’t cross-validated.

Explanation :

You can calculate the SSE by first making predictions and then computing the SSE:

The original Area tree was overfitting the data - it was uninterpretable. Area is not as useful as Murder - if it was, it would have been in the cross-validated tree. Cross-validation is not designed to improve the fit on the training data, but it won’t necessarily make it worse either. Cross-validation cannot guarantee improving the SSE on unseen data, although it often helps.