

<u>Help</u>



<u>Course</u> <u>Progress</u> <u>Dates</u> <u>Discussion</u> <u>Syllabus</u> <u>Schedule</u> <u>Files</u>

★ Course / Unit 4: Trees / Assignment 4

**(**)



### **State Data Revisted (OPTIONAL)**

 $\square$  Bookmark this page

# IMPORTANT NOTE: This problem is optional, and will not count towards your grade. We have created this problem to give you extra practice with the topics covered in this unit.

#### State data Revisited (OPTIONAL)

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:

data(state)
statedata = data.frame(state.x77)

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: <a href="mailto:statedataSimple.csv">statedataSimple.csv</a>. 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.

### Problem 1.1 - Linear Regression Models

0 points possible (ungraded)

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?

| 0.6922 |
|--------|
|        |

Explanation

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

RegModel =  $Im(Life.Exp \sim ., data=statedata)$ 

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

Submit

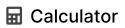
You have used 0 of 3 attempts



|   | ression Models  |
|---|---|
| D points possible (ungraded)<br>Calculate the sum of squared erro<br>actual life expectancies:  | ors (SSE) between the predicted life expectancies using this model and the  |
|   | Answer: 23.3  |
| Explanation To make predictions, type in your Predictions = predict(RegModel) where "RegModel" is the name of typing the following in your R con sum((statedata\$Life.Exp - Predic The SSE is 23.29714. Alternatively, you can use the foll sum(RegModel\$residuals^2)  | f your regression model. You can then compute the sum of squared errors by sole:<br>tions)^2)   |
| Submit You have used 0 of 3   | attempts  |
| Answers are displayed withing   | n the problem   |
|   |   |
| _   | model using just <i>Population, Murder, Frost, and HS.Grad</i> as independent del from the previous homework). What is the <b>adjusted</b> R-squared for this |
|   | Answer: 0.71  |
|   |   |
| ·<br>You can create this regression mo<br>RegModel2 = lm(Life.Exp ~ Popul<br>Then, if you type:<br>summary(RegModel2)   | odel by typing the following into your R console: ation + Murder + Frost + HS.Grad, data=statedata) bottom right of the output, and is 0.7126                 |
| ·<br>You can create this regression mo<br>RegModel2 = lm(Life.Exp ~ Popul<br>Then, if you type:<br>summary(RegModel2)   | ation + Murder + Frost + HS.Grad, data=statedata) bottom right of the output, and is 0.7126   |
| RegModel2 = Im(Life.Exp ~ Popul<br>Then, if you type:<br>summary(RegModel2)<br>The Adjusted R-squared is at the   | ation + Murder + Frost + HS.Grad, data=statedata)  bottom right of the output, and is 0.7126  attempts  |
| You can create this regression model and create this regression model and create this regression model and create the control of the control | ation + Murder + Frost + HS.Grad, data=statedata)  bottom right of the output, and is 0.7126  attempts  n the problem   |
| You can create this regression model RegModel 2 = Im(Life.Exp ~ Popul Then, if you type: summary(RegModel 2) The Adjusted R-squared is at the Submit You have used 0 of 3  Problem 1.4 - Linear Reg 0 points possible (ungraded)  | ation + Murder + Frost + HS.Grad, data=statedata)  bottom right of the output, and is 0.7126  attempts  n the problem   |

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

sum((statedata\$Life.Exp - Predictions2)^2).



| Submit   | You have used 0 of 3 attempts   |
|--|---|
| <b>3</b> Answe   | rs are displayed within the problem   |
| roblem   | 1.5 - Linear Regression Models  |
| points poss  | ible (ungraded) e following is correct?   |
|  |   |
|  | g different combinations of variables in linear regression is like trying different numbers of splits ree - this controls the complexity of the model.  |
| Using  | many variables in a linear regression is <b>always</b> better than using just a few.  |
| O The v  | ariables we removed were uncorrelated with <i>Life.Exp</i>  |
|  |   |
|  | of the model. This is similar to trying different numbers of splits in a tree, which is also controlling  |
| he second<br>djusted R-<br>he third an<br>ependent<br>0.59 and -0<br>e compute<br>or(stateda<br>or(stateda   | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. Is swer is incorrect because the variables we removed have non-zero correlations with the variable Life. Exp. Illiteracy and Area are negatively correlated with Life. Exp. with correlations of  |
| he second<br>djusted R-<br>he third an<br>ependent<br>0.59 and -0<br>e compute<br>or(stateda<br>or(stateda   | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. Is swer is incorrect because the variables we removed have non-zero correlations with the variable Life. Exp. Illiteracy and Area are negatively correlated with Life. Exp, with correlations of 0.11. Income is positively correlated with Life. Exp, with a correlation of 0.34. These correlations can do by typing the following into your R console: ta\$Life. Exp, statedata\$Income) ta\$Life. Exp, statedata\$Illiteracy)   |
| he second<br>djusted R-<br>he third an<br>ependent<br>0.59 and -0<br>e compute<br>or(stateda<br>or(stateda<br>or(stateda<br>Submit   | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. Is swer is incorrect because the variables we removed have non-zero correlations with the variable Life. Exp. Illiteracy and Area are negatively correlated with Life. Exp, with correlations of 0.11. Income is positively correlated with Life. Exp, with a correlation of 0.34. These correlations cand by typing the following into your R console: ta\$Life. Exp, statedata\$Income) ta\$Life. Exp, statedata\$Illiteracy) ta\$Life. Exp, statedata\$Area)   |
| he second djusted R-factorial divided in the second of the second of the computer of the conference of the second  | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. It is incorrect because the variables we removed have non-zero correlations with the variable Life.Exp. Illiteracy and Area are negatively correlated with Life.Exp, with correlations of 0.11. Income is positively correlated with Life.Exp, with a correlation of 0.34. These correlations cand by typing the following into your R console: ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Area)  You have used 0 of 1 attempt  |
| the second djusted R-fine third and lependent wood on the computer or (statedar or  | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. It is incorrect because the variables we removed have non-zero correlations with the variable Life.Exp. Illiteracy and Area are negatively correlated with Life.Exp, with correlations of D.11. Income is positively correlated with Life.Exp, with a correlation of 0.34. These correlations can do by typing the following into your R console:  ta\$Life.Exp, statedata\$Income)  ta\$Life.Exp, statedata\$Illiteracy)  ta\$Life.Exp, statedata\$Area)  You have used 0 of 1 attempt  2.1 - CART Models  ible (ungraded)  uild a CART model to predict Life.Exp using all of the other variables as independent variables income, liliteracy, Murder, HS.Grad, Frost, Area). We'll use the default minbucket parameter, so the minbucket argument. Remember that in this problem we are not as interested in predicting life to se for new observations as we are understanding how they relate to the other variables we have, all of the data to build our model. You shouldn't use the method="class" argument since this is a rece.  |
| the second djusted R-fine third and lependent woods and -0 oe compute for (statedator (sta | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. Isswer is incorrect because the variables we removed have non-zero correlations with the variable Life.Exp. Illiteracy and Area are negatively correlated with Life.Exp, with correlations of 0.11. Income is positively correlated with Life.Exp, with a correlation of 0.34. These correlations can d by typing the following into your R console: ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Area)  You have used 0 of 1 attempt  2.1 - CART Models  tible (ungraded)  uild a CART model to predict Life.Exp using all of the other variables as independent variables income, Illiteracy, Murder, HS.Grad, Frost, Area). We'll use the default minbucket parameter, so the minbucket argument. Remember that in this problem we are not as interested in predicting life in so for new observations as we are understanding how they relate to the other variables we have, all of the data to build our model. You shouldn't use the method="class" argument since this is a ree.  2. Which of these variables appear in the tree? Select all that apply. |
| the second djusted R-fine third and lependent wood on the computer or (statedar or  | answer is incorrect because as we see here, a model with fewer variables actually has a higher squared. If your accuracy is just as good, a model with fewer variables is almost always better. Isswer is incorrect because the variables we removed have non-zero correlations with the variable Life.Exp. Illiteracy and Area are negatively correlated with Life.Exp, with correlations of 0.11. Income is positively correlated with Life.Exp, with a correlation of 0.34. These correlations can d by typing the following into your R console: ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Income) ta\$Life.Exp, statedata\$Area)  You have used 0 of 1 attempt  2.1 - CART Models  tible (ungraded)  uild a CART model to predict Life.Exp using all of the other variables as independent variables income, Illiteracy, Murder, HS.Grad, Frost, Area). We'll use the default minbucket parameter, so the minbucket argument. Remember that in this problem we are not as interested in predicting life so for new observations as we are understanding how they relate to the other variables we have, all of the data to build our model. You shouldn't use the method="class" argument since this is a ree.  be Which of these variables appear in the tree? Select all that apply.    |

|  | ad   |
|--|--|
| Area   |  |
| CARTmodel<br>Be sure to lo<br>You can theo<br>orp(CARTmo | ate the tree in R by typing the following command: = rpart(Life.Exp ~ ., data=statedata) ad the "rpart" and "rpart.plot" packages with the library command if they are not already loaded. n plot the tree by typing: odel) that the only variable used in the tree is "Murder". |
| Submit   | You have used 0 of 3 attempts  |
| • Answei   | s are displayed within the problem   |
| Problem  | 2.2 - CART Models  |
| Jse the regr   | ble (ungraded) ession tree you just built to predict life expectancies (using the predict function), and calculate squared-errors (SSE) like you did for linear regression. What is the SSE?   |
|  | Answer: 29.0   |
| Γhen, you ca   | ART = predict(CARTmodel) an compute the sum of squared errors (SSE) by typing the following: ata\$Life.Exp - PredictionsCART)^2) 8.99848.  You have used 0 of 3 attempts   |
| 1 Answer   | s are displayed within the problem   |
| Problem  | 2.3 - CART Models  |
| The error is   | ble (ungraded) higher than for the linear regression models. One reason might be that we haven't made the tree Set the <i>minbucket</i> parameter to 5, and recreate the tree.   |
| Which varial   | ples appear in this new tree? Select all that apply.   |
| Popul  | ation  |
| ☐ Murde  | PT   |
| Frost  |  |
|  |  |

| Area ✓  |   |
|---|---|
|   |   |
| CARTmodel2 = rpart(Life.Exp ~ ., c  | ucket value of 5 with the following command:<br>data=statedata, minbucket=5)<br>o(CARTmodel2), you can see that Murder, HS.Grad, and Area are all used in |
| Submit You have used 0 of 3 a   | uttempts  |
| Answers are displayed within  | the problem   |
| Problem 2.4 - CART Mode   | els   |
| 0 points possible (ungraded)<br>Do you think the default minbucke   | t parameter is smaller or larger than 5 based on the tree that was built?   |
| Smaller   |   |
| Larger  |   |
| •   |   |
| from splitting more before. So the  Submit You have used 0 of 1 a   | default minbucket parameter must be larger than 5.  |
| Answers are displayed within  | the problem   |
| Problem 2.5 - CART Mode   | els   |
| 0 points possible (ungraded)<br>What is the SSE of this tree?   |   |
|   | Answer: 23.6  |
| Explanation You can compute the SSE of this to PredictionsCART2 = predict(CART) and then computing the sum of the sum((statedata\$Life.Exp - Predicti The SSE is 23.64283 | model2) e squared differences between the actual values and the predicted values:   |
| This is much closer to the linear re fit of our model.  | gression model's error. By changing the parameters we have improved the   |
| Submit You have used 0 of 3 a   | uttempts  |

• Answers are displayed within the problem

### Problem 2.6 - CART Models

0 points possible (ungraded)

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?

#### Explanation

You can create this third tree by typing:

CARTmodel3 = rpart(Life.Exp ~ Area, data=statedata, minbucket=1)

Then to compute the SSE, first make predictions:

PredictionsCART3 = predict(CARTmodel3)

And then compute the sum of squared differences between the actual values and the predicted values: sum((statedata\$Life.Exp - PredictionsCART3)^2)

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.

Submit

You have used 0 of 3 attempts

Answers are displayed within the problem

#### Problem 2.7 - CART Models

0 points possible (ungraded)

This is the lowest error we have seen so far. What would be the best interpretation of this result?

| $\bigcirc$ | Trees are much better than linear regress | ion for this problem | because they can | capture nonlinearities |
|------------|---|----------------------|------------------|------------------------|
|            | that linear regression misses.            |                      |                  |                        |

| $\bigcirc$ | We can build almost perfect models given the right parameters, even if they violate our intuition | of |
|------------|---|----|
|            | what a good model should be.  |    |



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.

Submit

You have used 0 of 1 attempt

**1** Answers are displayed within the problem

0 points possible (ungraded)

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 generalizeable. 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.)

| bottom of the output.)  |   |
|---|---|
|   | Answer: 0.12  |
| library(caret) set.seed(111) Then, you can set up the cross-va fitControl = trainControl(method = cartGrid = expand.grid(.cp = seq(0) You can then use the train function | "cv", number = 10) 0.01, 0.5, 0.01) ) n to find the best value of cp by typing: method="rpart", trControl = fitControl, tuneGrid = cartGrid)  |
| Submit You have used 0 of 4 a   | attempts  |
| Answers are displayed within  | the problem   |
| Problem 3.2 - Cross-Valid   | lation  |
| variables, and the entire dataset "   | you found in the previous problem, all of the available independent statedata" as the training data. Then plot the tree. You'll notice that this is see we created with the initial model. Interpret the tree: we predict the life rate is greater than or equal to |
|   | Answer: 6.6   |
|   |   |
| and is less than  |   |

### Explanation

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

CARTmodel4 = rpart(Life.Exp  $\sim$  ., data=statedata, cp=0.12)

Then, if you plot the tree using prp(CARTmodel4), 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).

Submit

You have used 0 of 4 attempts

**1** Answers are displayed within the problem

|  | Answer: 32.9  |
|--|---|
|  |   |
| Explanation<br>To compute                    | the SSE, first make predictions:  |
| •  | CART4 = predict(CARTmodel4)   |
|  | mpute the sum of squared differences between the actual values and the predicted values: ata\$Life.Exp - PredictionsCART4)^2) 32.86549  |
| Submit                                       | You have used 0 of 3 attempts   |
| <b>1</b> Answe                               | rs are displayed within the problem   |
| Problem                                      | 3.4 - Cross-Validation  |
| Recall the fi<br>validation) v<br>vou expect | rst tree (default parameters), second tree (minbucket = 5), and the third tree (selected with cross we made. Given what you have learned about cross-validation, which of the three models would to be better if we did use it for prediction on a test set? For this question, suppose we had actually few observations (states) in a test set, and we want to make predictions on those states. |
| The fi                                       | rst model   |
| The s  | econd model   |
| ○ The n                                      | nodel we just made with the "best" cp   |
| the model w                                  | e of cross-validation is to pick the tree that will perform the best on a test set. So we would expect<br>re made with the "best" cp to perform best on a test set.   |
| Submit                                       | You have used 0 of 1 attempt  |
| Answe  | rs are displayed within the problem   |
| Problem                                      | 3.5 - Cross-Validation  |
| At the end o                                 | ble (ungraded) If Problem 2 we made a very complex tree using just Area. Use <i>train</i> with the same parameters as ust using Area as an independent variable to find the best cp value (set the seed to 111 first). Then tree using just Area and this value of cp.  |
| How many s                                   | plits does the tree have?   |
|  | Answer: 4   |
|  |   |
| Explanation<br>To find the b                 | pest value of cp when using only Area, use the following command:   |

 $train(Life. Exp \sim Area, \ data = statedata, \ method = "rpart", \ trControl = fitControl, \ tuneGrid = cartGrid \ )$ 

Then, build a new CART tree by typina:

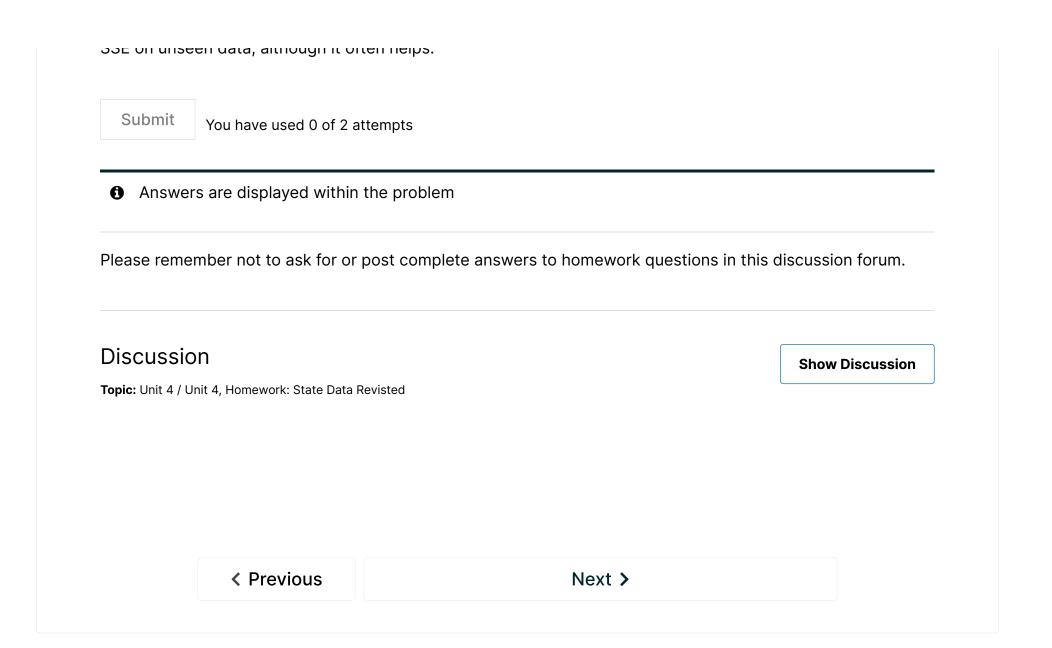
| Submit You have used 0   | of 4 attempts  |
|--|--|
| • Answers are displayed w  | ithin the problem  |
| Problem 3.6 - Cross-V  | alidation and the second secon |
| D points possible (ungraded)<br>The lower left leaf (or bucket)<br>correspond to states with are   | corresponds to the lowest predicted Life.Exp of 70. Observations in this leaf a greater than or equal to   |
|  | Answer: 9579   |
| and area less than   |  |
|  | Answer: 51000  |
| Explanation<br>To get to this leaf, we go thro<br>Area less than 62,000<br>Area greater than or equal to   |  |
| To get to this leaf, we go thro<br>Area less than 62,000<br>Area greater than or equal to<br>Area less than 51,000   | 9,579 mposed of states that have an area greater than 9,579 and less than 51,000.  |
| To get to this leaf, we go through the Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is co   | 9,579 mposed of states that have an area greater than 9,579 and less than 51,000. of 4 attempts  |
| Fo get to this leaf, we go through Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is consumpted Submit  You have used 0 to the Area less than 51,000  Submit  You have used 0 to the Area less than 51,000  | 9,579 mposed of states that have an area greater than 9,579 and less than 51,000. of 4 attempts rithin the problem   |
| To get to this leaf, we go throw Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is co  Submit  You have used 0  Problem 3.7 - Cross-V  Depoints possible (ungraded) We have simplified the previous   | 9,579 mposed of states that have an area greater than 9,579 and less than 51,000. of 4 attempts rithin the problem   |
| To get to this leaf, we go throw Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is consumption of the second of the s | mposed of states that have an area greater than 9,579 and less than 51,000.  of 4 attempts  of thin the problem  falidation  us "Area tree" considerably by using cross-validation. Calculate the SSE of the   |
| To get to this leaf, we go throw Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is consumpted by Submit  You have used 0  Problem 3.7 - Cross-V  Depoints possible (ungraded) We have simplified the previous cross-validated "Area tree", and the best model in this value.  | mposed of states that have an area greater than 9,579 and less than 51,000.  of 4 attempts  dithin the problem  alidation  us "Area tree" considerably by using cross-validation. Calculate the SSE of the ad select all of the following correct statements that apply:   |
| To get to this leaf, we go throw Area less than 62,000 Area greater than or equal to Area less than 51,000 This means that this leaf is consumpted by Submit  You have used 0  The have simplified the previous cross-validated "Area tree", and the best model in this way.  The Area variable is not the previous cross-validated and the previous cross-validated area tree.  | mposed of states that have an area greater than 9,579 and less than 51,000.  of 4 attempts  dithin the problem  alidation  us "Area tree" considerably by using cross-validation. Calculate the SSE of the nd select all of the following correct statements that apply:  whole question is the first "Area tree" because it had the lowest SSE.  as predictive as Murder rate.  |

PredictionsCART5 = predict(CARTmodel5)

sum((statedata\$Life.Exp - PredictionsCART5)^2)

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

Calculator



© All Rights Reserved



# edX

**About** 

<u>Affiliates</u>

edX for Business

Open edX

Careers

News

# Legal

Terms of Service & Honor Code

Privacy Policy

Accessibility Policy

Trademark Policy

<u>Sitemap</u>

Cookie Policy

**Your Privacy Choices** 



# **Connect**

<u>Idea Hub</u>

Contact Us

Help Center

<u>Security</u>

Media Kit















© 2024 edX LLC. All rights reserved. 深圳市恒宇博科技有限公司 <u>粤ICP备17044299号-2</u>