**edX** (https://www.edx.org)

MITx: 15.071x The Analytics Edge        ggsatz (/dashboard)   ▼

Courseware (/courses/MITx/15.071x/1T2014/courseware)      Course Info (/courses/MITx/15.071x/1T2014/info)

Discussion (/courses/MITx/15.071x/1T2014/discussion/forum)      Progress (/courses/MITx/15.071x/1T2014/progress)

Syllabus (/courses/MITx/15.071x/1T2014/4264e68418f34d839cf0b33a5da644b2/)

Schedule (/courses/MITx/15.071x/1T2014/2891f8bf120945b9aa12e6601739c3e6/)

Help

## STATE DATA

We often take data for granted. However, one of the hardest parts about analyzing a problem you're interested in can be to find good data to answer the questions you want to ask. As you're learning R, though, there are many datasets that R has built in that you can take advantage of.

In this problem, we will be examining the "state" dataset, which has data from the 1970s on all fifty US states. For each state, the dataset includes 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.

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

data(state)

statedata = cbind(data.frame(state.x77), state.abb, state.area, state.center,  state.division, state.name, state.region)

Inspect the data set using the command: str(statedata)

For more information about this data set, type ?state in the R console.

## PROBLEM 1.1 - DATA EXPLORATION  (1/1 point)

We begin by exploring the data by examining the latitude and longitude of each state. Plot all of the states' centers with latitude on the y axis (the "y" variable in our dataset) and longitude on the x axis (the "x" variable in our dataset). The shape of the plot should be the familiar outline of the United States! Note that Alaska and Hawaii have had their coordinates adjusted to appear just off of the west coast.

In the R command you used to generate this plot, which variable name did you use as the first argument?

- ○ statedata$y
- ◉ statedata$x   ✔
- ○ I used a different variable name.

**EXPLANATION**

To generate the described plot, you should type plot(statedata$x, statedata$y) in your R console. The first variable here is statedata$x.

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 1.2 - DATA EXPLORATION  (1/1 point)

Using the tapply command, determine which region of the US (West, North Central, South, or Northeast) has the highest average high school graduation rate of all the states in the region:

- ⦿ West ✔
- ○ North Central
- ○ South
- ○ Northeast

**EXPLANATION**

You can find the average high school graduation rate of all states in each of the regions by typing the following command in your R console:

tapply(statedata$HS.Grad, state.region, mean)

The highest value is for the West region.

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 1.3 - DATA EXPLORATION  (1/1 point)

Now, make a boxplot of the murder rate by region (for more information about creating boxplots in R, type ?boxplot in your console).

Which region has the highest median murder rate?

- ○ Northeast
- ⦿ South ✔
- ○ North Central
- ○ West

**EXPLANATION**

To generate the boxplot, you should type boxplot(statedata$Murder ~ statedata$state.region) in your R console. You can see that the region with the highest median murder rate (the one with the highest solid line in the box) is the South.

Help

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 1.4 - DATA EXPLORATION  (1/1 point)

You should see that there is an outlier in the Northeast region of the boxplot you just generated. Which state does this correspond to? (Hint: There are many ways to find the answer to this question, but one way is to use the subset command to only look at the Northeast data.)

- ○ Delaware
- ○ Rhode Island
- ○ Maine
- ● New York ✔

**EXPLANATION**

The correct answer is New York. If you first use the subset command:

NortheastData = subset(statedata, state.region == "Northeast")

Then look at NortheastData$Murder together with NortheastData$state.abb to identify the outlier.

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 2.1 - PREDICTING LIFE EXPECTANCY - AN INITIAL MODEL  (1/1 point)

We would like to build a model to predict life expectancy by state using the state statistics we have in our dataset.

Build the model with all potential variables included (Population, Income, Illiteracy, Murder, HS.Grad, Frost, and Area). Note that you should use the variable "Area" in your model, NOT the variable "state.area".

What is coefficient for income?

-2.180e-05

$$-2.180 \times 10^{-05}$$

**Answer:** -0.0000218

**EXPLANATION**

You can build the linear regression model with the following command:

LinReg = lm(Life.Exp ~ Population + Income + Illiteracy + Murder + HS.Grad + Frost + Area, data=statedata)

Then, to find the coefficient for income, you can look at the summary of the regression with summary(LinReg).

Hide Answer    *You have used 1 of 3 submissions*

## PROBLEM 2.2 - PREDICTING LIFE EXPECTANCY - AN INITIAL MODEL  (1/1 point)

Call the coefficient for income x (the answer to Problem 2.1). What is the interpretation of the coefficient x?

- ○ For a one unit increase in income, predicted life expectancy increases by |x|
- ● For a one unit increase in income, predicted life expectancy decreases by |x|  ✔
- ○ For a one unit increase in predicted life expectancy, income decreases by |x|
- ○ For a one unit increase in predicted life expectancy, income increases by |x|

**EXPLANATION**

If we increase income by one unit, then our model's prediction will increase by the coefficient of income, x. Because x is negative, this is the same as predicted life expectancy decreasing by |x|.

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 2.3 - PREDICTING LIFE EXPECTANCY - AN INITIAL MODEL  (1/1 point)

Now plot a graph of life expectancy vs. income using the command:

plot(statedata$Income, statedata$Life.Exp)

Visually observe the plot. What appears to be the relationship?

- ● Life expectancy is somewhat positively correlated with income.  ✔
- ○ Life expectancy is somewhat negatively correlated with income.
- ○ Life expectancy is not correlated with income.

**EXPLANATION**

Although the point in the lower right hand corner of the plot appears to be an outlier, we observe a positive linear relationship in the plot.

Hide Answer    *You have used 1 of 1 submissions*

## PROBLEM 2.4 - PREDICTING LIFE EXPECTANCY - AN INITIAL MODEL  (1 point possible)

The model we built does not display the relationship we saw from the plot of life expectancy vs. income. Which of the

following explanations seems the most reasonable?

    &#9675; Income is not related to life expectancy.

    &#9673; Multicollinearity

**EXPLANATION**

Although income is an insignificant variable in the model, this does not mean that there is no association between income and life expectancy. However, in the presence of all of the other variables, income does not add statistically significant explanatory power to the model. This means that multicollinearity is probably the issue.

Hide Answer  *You have used 1 of 1 submissions*

## PROBLEM 3.1 - PREDICTING LIFE EXPECTANCY - REFINING THE MODEL AND ANALYZING PREDICTIONS

 (1/1 point)

Recall that we discussed the principle of simplicity: that is, a model with fewer variables is preferable to a model with many unnnecessary variables. Experiment with removing independent variables from the original model. Remember to use the significance of the coefficients to decide which variables to remove (remove the one with the largest "p-value" first, or the one with the "t value" closest to zero), and to remove them one at a time (this is called "backwards variable selection"). This is important due to multicollinearity issues - removing one insignificant variable may make another previously insignificant variable become significant.

You should be able to find a good model with only 4 independent variables, instead of the original 7. Which variables does this model contain?

    &#9675; Income, HS.Grad, Frost, Murder

    &#9675; HS.Grad, Population, Income, Frost

    &#9675; Frost, Murder, HS.Grad, Illiteracy

    &#9673; Population, Murder, Frost, HS.Grad ✔

**EXPLANATION**

We would eliminate the variable "Area" first (since it has the highest p-value, or probability, with a value of 0.9649), by adjusting our lm command to the following:

LinReg = lm(Life.Exp ~ Population + Income + Illiteracy + Murder + HS.Grad + Frost + Area, data=statedata)

Looking at summary(LinReg) now, we would choose to eliminate "Illiteracy" since it now has the highest p-value of 0.9340, using the following command:

LinReg = lm(Life.Exp ~ Population + Income + Murder + HS.Grad + Frost + Area, data=statedata)

Looking at summary(LinReg) again, we would next choose to eliminate "Income", since it has a p-value of 0.9153. This gives the following four variable model:

LinReg = lm(Life.Exp ~ Population + Murder + HS.Grad + Frost, data=statedata)

This model with 4 variables is a good model. However, we can see that the variable "Population" is not quite significant. In practice, it would be up to you whether or not to keep the variable "Population" or eliminate it for a 3-variable model. Population does not add much statistical significance in the presence of murder, high school graduation rate, and frost days. However, for the remainder of this question, we will analyze the 4-variable model.

**Hide Answer**      *You have used 1 of 1 submissions*

## PROBLEM 3.2 - PREDICTING LIFE EXPECTANCY - REFINING THE MODEL AND ANALYZING PREDICTIONS

(1 point possible)

Removing insignificant variables changes the Multiple R-squared value of the model. By looking at the summary output for both the initial model (all independent variables) and the simplified model (only 4 independent variables) and using what you learned in class, which of the following correctly explains the change in the Multiple R-squared value?

- ⦿ We expect the "Multiple R-squared" value of the simplified model to be slightly worse than that of the initial model. It can't be better than the "Multiple R-squared" value of the initial model.
- ○ We expect the "Multiple R-squared" value of the simplified model to be slightly better than that of the initial model. It can't be worse than the "Multiple R-squared" value of the initial model.
- ○ We expect the "Multiple R-squared" of the simplified model to be about the same as the intial model (we have no way of knowing if it will be slightly worse or slightly better than the Multiple R-squared of the intial model).

**EXPLANATION**

When we remove insignificant variables, the "Multiple R-squared" will always be worse, but only slightly worse. This is due to the nature of a linear regression model. It is always possible for the regression model to make a coefficient zero, which would be the same as removing the variable from the model. The fact that the coefficient is not zero in the intial model means it must be helping the R-squared value, even if it is only a very small improvement. So when we force the variable to be removed, it will decrease the R-squared a little bit. However, this small decrease is worth it to have a simpler model.

On the contrary, when we remove insignificant variables, the "Adjusted R-squred" will frequently be better. This value accounts for the complexity of the model, and thus tends to increase as insignificant variables are removed, and decrease as insignificant variables are added.

**Hide Answer**      *You have used 2 of 2 submissions*

## PROBLEM 3.3 - PREDICTING LIFE EXPECTANCY - REFINING THE MODEL AND ANALYZING PREDICTIONS

(1/2 points)

Using the simplified 4 variable model that we created, we'll now take a look at how our predictions compare to the actual values.

Take a look at the vector of predictions by using the predict function (since we are just looking at predictions on the training set, you don't need to pass a "newdata" argument to the predict function).

Which state do we predict to have the lowest life expectancy? (Hint: use the sort function)

- ○ South Carolina
- ○ Mississippi
- ◉ Alabama
- ○ Georgia

---

**EXPLANATION**

If your simplified 4-variable model is called "LinReg", you can answer this question by typing

sort(predict(LinReg))

in your R console. The first state listed has the lowest predicted life expectancy, which is Alabama.

---

Which state actually has the lowest life expectancy? (Hint: use the which.min function)

- ◉ South Carolina ✔
- ○ Mississippi
- ○ Alabama
- ○ Georgia

---

**EXPLANATION**

You can find the row number of the state with the lowest life expectancy by typing which.min(statedata$Life.Exp) into your R console. This returns 40. The 40th state name in the vector statedata$state.name is South Carolina.

---

| Hide Answer |    *You have used 1 of 1 submissions*

---

## PROBLEM 3.4 - PREDICTING LIFE EXPECTANCY - REFINING THE MODEL AND ANALYZING PREDICTIONS

(2/2 points)

Which state do we predict to have the highest life expectancy?

    ○ Massachusetts

    ○ Maine

    ◉ Washington    ✔

    ○ Hawaii

Which state actually has the highest life expectancy?

    ○ Massachusetts

    ○ Maine

    ○ Washington

    ◉ Hawaii    ✔

| Show Answer | *You have used 1 of 1 submissions* |
|---|---|

## PROBLEM 3.5 - PREDICTING LIFE EXPECTANCY - REFINING THE MODEL AND ANALYZING PREDICTIONS

 (1/2 points)

Take a look at the vector of residuals (the difference between the predicted and actual values).

For which state do we make the smallest absolute error?

    ○ Maine

    ○ Florida

    ◉ Indiana

    ○ Illinois

---

**EXPLANATION**

You can look at the sorted list of absolute errors by typing

sort(abs(model$residuals))

into your R console (where "model" is the name of your model). Alternatively, you can compute the residuals manually by typing

sort(abs(statedata$Life.Exp - predict(model)))

in your R console. The smallest absolute error is for Indiana.

---

For which state do we make the largest absolute error?

⊙ Hawaii　✔

○ Maine

○ Texas

○ South Carolina

**Help**

**EXPLANATION**

You can look at the sorted list of absolute errors by typing

sort(abs(model$residuals))

into your R console (where "model" is the name of your model). Alternatively, you can compute the residuals manually by typing

sort(abs(statedata$Life.Exp - predict(model)))

in your R console. The largest absolute error is for Hawaii.

| Hide Answer |   *You have used 1 of 1 submissions*

Please remember not to ask for or post complete answers to homework questions in this discussion forum.

Show Discussion　　　　　　　　　　　　　　✎ **New Post**

(http://www.meetup.com/edX-Global-Community/)

(http://www.facebook.com/EdxOnline)

(https://twitter.com/edXOnline)

(https://plus.google.com/108235383044095082735/posts)

Help