Unit 2: Linear Regression > Assignment 2 > Reading Test Scores



#### MITx: 15.071x The Analytics Edge



■ Bookmark

- Unit 1: An Introduction to Analytics
- EntranceSurvey
- ▼ Unit 2: Linear Regression

#### Welcome to Unit 2

The Statistical
Sommelier: An
Introduction to
Linear Regression
Lecture Sequence
Quick Questions

Moneyball: The Power of Sports Analytics

Lecture Sequence
Quick Questions

Ø.

Playing Moneyball in the NBA (Recitation)

#### **Assignment 2**

Homework due May 03, 2016 at 00:00 UT

Unit 3: Logistic Regression

#### **READING TEST SCORES**

The Programme for International Student Assessment (PISA) is a test given every three years to 15-year-old students from around the world to evaluate their performance in mathematics, reading, and science. This test provides a quantitative way to compare the performance of students from different parts of the world. In this homework assignment, we will predict the reading scores of students from the United States of America on the 2009 PISA exam.

The datasets <u>pisa2009train.csv</u> and <u>pisa2009test.csv</u> contain information about the demographics and schools for American students taking the exam, derived from <u>2009 PISA Public-Use Data Files</u> distributed by the United States National Center for Education Statistics (NCES). While the datasets are not supposed to contain identifying information about students taking the test, by using the data you are bound by the <u>NCES data use agreement</u>, which prohibits any attempt to determine the identity of any student in the datasets.

Each row in the datasets pisa2009train.csv and pisa2009test.csv represents one student taking the exam. The datasets have the following variables:

**grade:** The grade in school of the student (most 15-year-olds in America are in 10th grade)

male: Whether the student is male (1/0)

raceeth: The race/ethnicity composite of the student

preschool: Whether the student attended preschool (1/0)

**expectBachelors:** Whether the student expects to obtain a bachelor's degree (1/0)

motherHS: Whether the student's mother completed high school (1/0)

**motherBachelors:** Whether the student's mother obtained a bachelor's degree (1/0)

**motherWork:** Whether the student's mother has part-time or full-time work (1/0)

**fatherHS:** Whether the student's father completed high school (1/0)

**fatherBachelors:** Whether the student's father obtained a bachelor's degree (1/0)

**fatherWork:** Whether the student's father has part-time or full-time work (1/0)

**selfBornUS:** Whether the student was born in the United States of America (1/0)

**motherBornUS:** Whether the student's mother was born in the United States of America (1/0)

**fatherBornUS:** Whether the student's father was born in the United States of America (1/0)

englishAtHome: Whether the student speaks English at home (1/0)

**computerForSchoolwork:** Whether the student has access to a computer for schoolwork (1/0)

**read30MinsADay:** Whether the student reads for pleasure for 30 minutes/day (1/0)

**minutesPerWeekEnglish:** The number of minutes per week the student spend in English class

**studentsInEnglish:** The number of students in this student's English class at school

schoolHasLibrary: Whether this student's school has a library (1/0)

**publicSchool:** Whether this student attends a public school (1/0)

urban: Whether this student's school is in an urban area (1/0)

**schoolSize:** The number of students in this student's school

readingScore: The student's reading score, on a 1000-point scale

### Problem 1.1 - Dataset size

(1/1 point)

Load the training and testing sets using the read.csv() function, and save them as variables with the names pisaTrain and pisaTest.

How many students are there in the training set?

3663



You have used 2 of 3 submissions

# Problem 1.2 - Summarizing the dataset

(2/2 points)

Using tapply() on pisaTrain, what is the average reading test score of males?

483.532478632479



Of females?

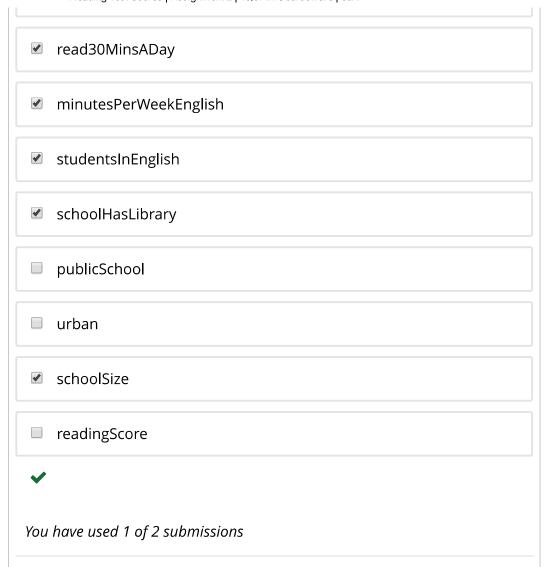
512.94063093244



You have used 1 of 3 submissions

# Problem 1.3 - Locating missing values

(1/1 point) Which variables are missing data in at least one observation in the training set? Select all that apply.		
	grade	
	male	
•	raceeth	
•	preschool	
•	expectBachelors	
•	motherHS	
•	motherBachelors	
•	motherWork	
•	fatherHS	
•	fatherBachelors	
•	fatherWork	
•	selfBornUS	
•	motherBornUS	
•	fatherBornUS	
•	englishAtHome	
•	computerForSchoolwork	



# Problem 1.4 - Removing missing values

(2/2 points)

Linear regression discards observations with missing data, so we will remove all such observations from the training and testing sets. Later in the course, we will learn about imputation, which deals with missing data by filling in missing values with plausible information.

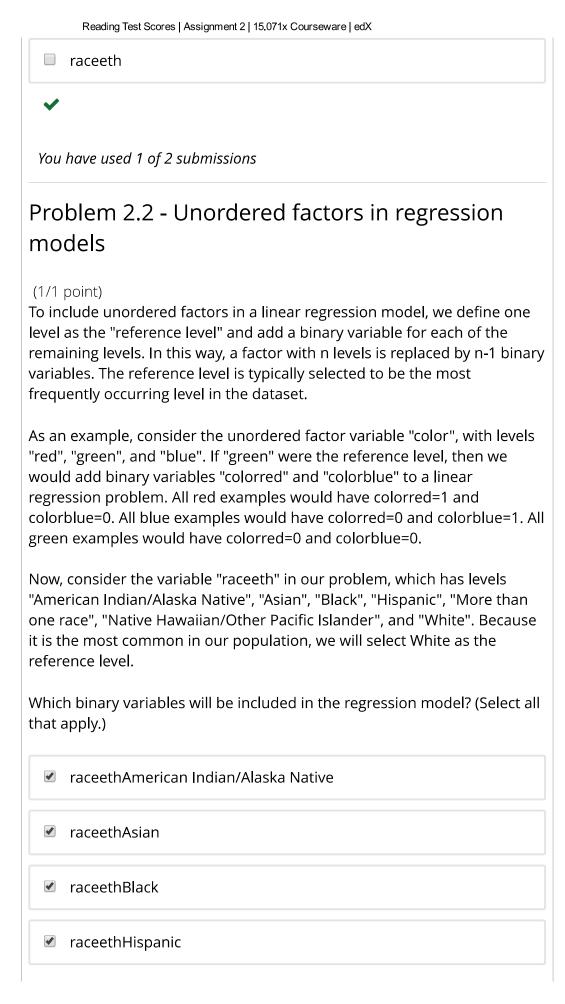
Type the following commands into your R console to remove observations with any missing value from pisaTrain and pisaTest:

pisaTrain = na.omit(pisaTrain)

pisaTest = na.omit(pisaTest)

How many observations are now in the training set?

24	14	<b>✓</b>
How	many observations are	e now in the testing set?
990	)	✓
You	have used 1 of 3 submis	ssions
Pro	blem 2.1 - Facto	or variables
Factor "Reg This betw level "sma	ion" variable in the WH is an unordered factor een the levels. An orde s (an example would be ill").	es that take on a discrete set of values, like the O dataset from the second lecture of Unit 1. because there isn't any natural ordering red factor has a natural ordering between the e the classifications "large," "medium," and bles is an unordered factor with at least 3
	grade	,
	male	
•	raceeth	
	h of the following varia ct all that apply.)	bles is an ordered factor with at least 3 levels?
•	grade	
	male	



✓ raceethMore than one race		
✓ raceethNative Hawaiian/Other Pacific Islander		
raceethWhite		
<b>✓</b>		
You have used 2 of 2 submissions		
Problem 2.3 - Example unordered factors		
(2/2 points)  Consider again adding our unordered factor race to the regression model with reference level "White".		
For a student who is Asian, which binary variables would be set to 0? All remaining variables will be set to 1. (Select all that apply.)		
☑ raceethAmerican Indian/Alaska Native		
□ raceethAsian		
✓ raceethBlack		
✓ raceethHispanic		
☑ raceethMore than one race		
☑ raceethNative Hawaiian/Other Pacific Islander		
<b>✓</b>		
For a student who is white, which binary variables would be set to 0? All remaining variables will be set to 1. (Select all that apply.)		
raceethAmerican Indian/Alaska Native		

- ✓ raceethAsian
   ✓ raceethBlack
   ✓ raceethHispanic
   ✓ raceethMore than one race
   ✓ raceethNative Hawaiian/Other Pacific Islander
  - **V**

You have used 1 of 2 submissions

# Problem 3.1 - Building a model

(2/2 points)

Because the race variable takes on text values, it was loaded as a factor variable when we read in the dataset with read.csv() -- you can see this when you run str(pisaTrain) or str(pisaTest). However, by default R selects the first level alphabetically ("American Indian/Alaska Native") as the reference level of our factor instead of the most common level ("White"). Set the reference level of the factor by typing the following two lines in your R console:

pisaTrain\$raceeth = relevel(pisaTrain\$raceeth, "White")

pisaTest\$raceeth = relevel(pisaTest\$raceeth, "White")

Now, build a linear regression model (call it lmScore) using the training set to predict readingScore using all the remaining variables.

It would be time-consuming to type all the variables, but R provides the shorthand notation "readingScore ~ ." to mean "predict readingScore using all the other variables in the data frame." The period is used to replace listing out all of the independent variables. As an example, if your dependent variable is called "Y", your independent variables are called "X1", "X2", and "X3", and your training data set is called "Train", instead of the regular notation:

LinReg =  $Im(Y \sim X1 + X2 + X3, data = Train)$ 

You would use the following command to build your model:

 $LinReg = Im(Y \sim ., data = Train)$ 

What is the Multiple R-squared value of ImScore on the training set?

0.3251



Note that this R-squared is lower than the ones for the models we saw in the lectures and recitation. This does not necessarily imply that the model is of poor quality. More often than not, it simply means that the prediction problem at hand (predicting a student's test score based on demographic and school-related variables) is more difficult than other prediction problems (like predicting a team's number of wins from their runs scored and allowed, or predicting the quality of wine from weather conditions).

You have used 1 of 5 submissions

# Problem 3.2 - Computing the root-mean squared error of the model

(1 point possible)

What is the training-set root-mean squared error (RMSE) of ImScore?

374.196554323944



**Answer:** 73.36555

#### **EXPLANATION**

The training-set RMSE can be computed by first computing the SSE:

SSE = sum(ImScore\$residuals^2)

and then dividing by the number of observations and taking the square root:

RMSE = sqrt(SSE / nrow(pisaTrain))

A alternative way of getting this answer would be with the following command:

sqrt(mean(lmScore\$residuals^2)).

You have used 3 of 3 submissions

# Problem 3.3 - Comparing predictions for similar students

(1 point possible)

Consider two students A and B. They have all variable values the same, except that student A is in grade 11 and student B is in grade 9. What is the predicted reading score of student A minus the predicted reading score of student B?

- -59.09
- -29.54
- 0 X
- 29.54
- The difference cannot be determined without more information about the two students

#### **EXPLANATION**

The coefficient 29.54 on grade is the difference in reading score between two students who are identical other than having a difference in grade of 1. Because A and B have a difference in grade of 2, the model predicts that student A has a reading score that is 2\*29.54 larger.

You have used 2 of 2 submissions

### Problem 3.4 - Interpreting model coefficients

(1 point possible)

What is the meaning of the coefficient associated with variable raceethAsian?

- Predicted average reading score of an Asian student
- O Difference between the average reading score of an Asian student and the average reading score of a white student
- Difference between the average reading score of an Asian studentand the average reading score of all the students in the dataset
- Predicted difference in the reading score between an Asian student and a white student who is otherwise identical ✓

#### **EXPLANATION**

The only difference between an Asian student and white student with otherwise identical variables is that the former has raceethAsian=1 and the latter has raceethAsian=0. The predicted reading score for these two students will differ by the coefficient on the variable raceethAsian.

You have used 1 of 1 submissions

# Problem 3.5 - Identifying variables lacking statistical significance

(1/1 point)

remo	d on the significance codes, which variables are candidates for oval from the model? Select all that apply. (We'll assume that the r variable raceeth should only be removed if none of its levels are ficant.)
	grade
	male
	raceeth
<b>✓</b>	preschool
	expectBachelors
•	motherHS
	motherBachelors
•	motherWork
•	fatherHS
	fatherBachelors
•	fatherWork
•	selfBornUS
•	motherBornUS
•	fatherBornUS
•	englishAtHome

□ computerForSchoolwork		
read30MinsADay		
minutesPerWeekEnglish		
✓ studentsInEnglish		
✓ schoolHasLibrary		
publicSchool		
<b>☑</b> urban		
schoolSize		
<b>✓</b>		
You have used 1 of 2 submissions		
Problem 4.1 - Predicting on unseen data		
(2/2 points) Using the "predict" function and supplying the "newdata" argument, use the ImScore model to predict the reading scores of students in pisaTest. Call this vector of predictions "predTest". Do not change the variables in the model (for example, do not remove variables that we found were not significant in the previous part of this problem). Use the summary function to describe the test set predictions.		
What is the range between the maximum and minimum predicted reading score on the test set?		
284.5		
You have used 2 of 5 submissions		

### Problem 4.2 - Test set SSE and RMSE

(2/2 points)

What is the sum of squared errors (SSE) of ImScore on the testing set?

5762082.3711444



What is the root-mean squared error (RMSE) of ImScore on the testing set?

76.290793831092



You have used 1 of 5 submissions

# Problem 4.3 - Baseline prediction and test-set SSE

(2/2 points)

What is the predicted test score used in the baseline model? Remember to compute this value using the training set and not the test set.

517.962887323944



What is the sum of squared errors of the baseline model on the testing set? HINT: We call the sum of squared errors for the baseline model the total sum of squares (SST).

7802354.07761384



You have used 1 of 5 submissions

# Problem 4.4 - Test-set R-squared

(1/1 point)

What is the test-set R-squared value of ImScore?

0.261494375437702



You have used 1 of 5 submissions

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

© All Rights Reserved



© edX Inc. All rights reserved except where noted. EdX, Open edX and the edX and Open EdX logos are registered trademarks or trademarks of edX Inc.

















