

R Notebook

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.4.2
```

```
## -- Attaching packages -----
```

```
## v ggplot2 3.1.0      v purrr   0.2.5
## v tibble  2.0.1      v dplyr  0.7.8
## v tidyr   0.8.0      v stringr 1.3.1
## v readr   1.1.1      v forcats 0.3.0
```

```
## Warning: package 'ggplot2' was built under R version 3.4.4
```

```
## Warning: package 'tibble' was built under R version 3.4.4
```

```
## Warning: package 'tidyr' was built under R version 3.4.3
```

```
## Warning: package 'purrr' was built under R version 3.4.4
```

```
## Warning: package 'dplyr' was built under R version 3.4.4
```

```
## Warning: package 'forcats' was built under R version 3.4.3
```

```
## -- Conflicts -----
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
library(skimr)
```

```
## Warning: package 'skimr' was built under R version 3.4.4
```

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 **pisa2009train.csv** and **pisa2009test.csv** contain information about the demographics and schools for American students taking the exam, derived from 2009 PISA Public-Use Data Files 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 NCES data use agreement, 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

1.1 Dataset size

```
pisa_train = read.csv('pisa2009train.csv')
pisa_test = read.csv('pisa2009test.csv')

str(pisa_train)
```

```
## 'data.frame':   3663 obs. of  24 variables:
## $ grade          : int  11 11 9 10 10 10 10 10 9 10 ...
## $ male           : int  1 1 1 0 1 1 0 0 0 1 ...
## $ raceeth        : Factor w/ 7 levels "American Indian/Alaska Native",...: NA 7 7 3 4 3 2 7 7 5
## $ preschool      : int  NA 0 1 1 1 1 0 1 1 1 ...
## $ expectBachelors: int  0 0 1 1 0 1 1 1 0 1 ...
## $ motherHS       : int  NA 1 1 0 1 NA 1 1 1 1 ...
## $ motherBachelors: int  NA 1 1 0 0 NA 0 0 NA 1 ...
## $ motherWork     : int  1 1 1 1 1 1 1 0 1 1 ...
## $ fatherHS       : int  NA 1 1 1 1 1 NA 1 0 0 ...
## $ fatherBachelors: int  NA 0 NA 0 0 0 NA 0 NA 0 ...
## $ fatherWork     : int  1 1 1 1 0 1 NA 1 1 1 ...
## $ selfBornUS     : int  1 1 1 1 1 1 0 1 1 1 ...
## $ motherBornUS   : int  0 1 1 1 1 1 1 1 1 1 ...
## $ fatherBornUS   : int  0 1 1 1 0 1 NA 1 1 1 ...
## $ englishAtHome  : int  0 1 1 1 1 1 1 1 1 1 ...
## $ computerForSchoolwork: int  1 1 1 1 1 1 1 1 1 1 ...
## $ read30MinsADay : int  0 1 0 1 1 0 0 1 0 0 ...
## $ minutesPerWeekEnglish: int  225 450 250 200 250 300 250 300 378 294 ...
## $ studentsInEnglish : int  NA 25 28 23 35 20 28 30 20 24 ...
## $ schoolHasLibrary : int  1 1 1 1 1 1 1 0 1 ...
```

```
## $ publicSchool      : int  1 1 1 1 1 1 1 1 1 1 ...
## $ urban             : int  1 0 0 1 1 0 1 0 1 0 ...
## $ schoolSize        : int  673 1173 1233 2640 1095 227 2080 1913 502 899 ...
## $ readingScore      : num  476 575 555 458 614 ...
```

1.2 Summarizing the dataset

Using `tapply()` on `pisaTrain`, what is the average reading test score of males (Ans: 483.5325)

```
tapply(pisa_train$readingScore, pisa_train$male == 1, mean)
```

```
##      FALSE      TRUE
## 512.9406 483.5325
```

1.3 Locating missing values

```
pisa_train %>%
  is.na() %>%
  sum()
```

```
## [1] 2950
```

```
summary(pisa_train)
```

```
##      grade      male      raceeth
## Min.   : 8.00   Min.   :0.0000   White      :2015
## 1st Qu.:10.00   1st Qu.:0.0000   Hispanic    : 834
## Median :10.00   Median :1.0000   Black       : 444
## Mean   :10.09   Mean   :0.5111   Asian       : 143
## 3rd Qu.:10.00   3rd Qu.:1.0000   More than one race: 124
## Max.   :12.00   Max.   :1.0000   (Other)     :  68
##                                     NA's       :  35
##      preschool    expectBachelors    motherHS    motherBachelors
## Min.   :0.0000   Min.   :0.0000   Min.   :0.00   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:1.0000   1st Qu.:1.00   1st Qu.:0.0000
## Median :1.0000   Median :1.0000   Median :1.00   Median :0.0000
## Mean   :0.7228   Mean   :0.7859   Mean   :0.88   Mean   :0.3481
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.00   3rd Qu.:1.0000
## Max.   :1.0000   Max.   :1.0000   Max.   :1.00   Max.   :1.0000
## NA's    :56     NA's    :62     NA's    :97     NA's    :397
##      motherWork    fatherHS    fatherBachelors    fatherWork
## Min.   :0.0000   Min.   :0.0000   Min.   :0.0000   Min.   :0.0000
## 1st Qu.:0.0000   1st Qu.:1.0000   1st Qu.:0.0000   1st Qu.:1.0000
## Median :1.0000   Median :1.0000   Median :0.0000   Median :1.0000
## Mean   :0.7345   Mean   :0.8593   Mean   :0.3319   Mean   :0.8531
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000
## Max.   :1.0000   Max.   :1.0000   Max.   :1.0000   Max.   :1.0000
## NA's    :93     NA's    :245    NA's    :569    NA's    :233
##      selfBornUS    motherBornUS    fatherBornUS    englishAtHome
## Min.   :0.0000   Min.   :0.0000   Min.   :0.0000   Min.   :0.0000
## 1st Qu.:1.0000   1st Qu.:1.0000   1st Qu.:1.0000   1st Qu.:1.0000
## Median :1.0000   Median :1.0000   Median :1.0000   Median :1.0000
## Mean   :0.9313   Mean   :0.7725   Mean   :0.7668   Mean   :0.8717
## 3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000   3rd Qu.:1.0000
```

```
## Max. :1.0000 Max. :1.0000 Max. :1.0000 Max. :1.0000
## NA's :69 NA's :71 NA's :113 NA's :71
## computerForSchoolwork read30MinsADay minutesPerWeekEnglish
## Min. :0.0000 Min. :0.0000 Min. : 0.0
## 1st Qu.:1.0000 1st Qu.:0.0000 1st Qu.: 225.0
## Median :1.0000 Median :0.0000 Median : 250.0
## Mean :0.8994 Mean :0.2899 Mean : 266.2
## 3rd Qu.:1.0000 3rd Qu.:1.0000 3rd Qu.: 300.0
## Max. :1.0000 Max. :1.0000 Max. :2400.0
## NA's :65 NA's :34 NA's :186
## studentsInEnglish schoolHasLibrary publicSchool urban
## Min. : 1.0 Min. :0.0000 Min. :0.0000 Min. :0.0000
## 1st Qu.:20.0 1st Qu.:1.0000 1st Qu.:1.0000 1st Qu.:0.0000
## Median :25.0 Median :1.0000 Median :1.0000 Median :0.0000
## Mean :24.5 Mean :0.9676 Mean :0.9339 Mean :0.3849
## 3rd Qu.:30.0 3rd Qu.:1.0000 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :75.0 Max. :1.0000 Max. :1.0000 Max. :1.0000
## NA's :249 NA's :143
## schoolSize readingScore
## Min. : 100 Min. :168.6
## 1st Qu.: 712 1st Qu.:431.7
## Median :1212 Median :499.7
## Mean :1369 Mean :497.9
## 3rd Qu.:1900 3rd Qu.:566.2
## Max. :6694 Max. :746.0
## NA's :162
```

1.4) Removing missing values

```
pisa_train = na.omit(pisa_train)
pisa_test = na.omit(pisa_test)
```

2.1) Factor variables

Factor variables are variables that take on a discrete set of values, like the “Region” variable in the WHO dataset from the second lecture of Unit 1. This is an unordered factor because there isn’t any natural ordering between the levels. An ordered factor has a natural ordering between the levels (an example would be the classifications “large,” “medium,” and “small”).

```
# These are unordered factors w/ at least 3 levels
table(pisa_train$male)
```

```
##
##      0      1
## 1204 1210
```

```
table(pisa_train$raceeth)
```

```
##
##      American Indian/Alaska Native
##                      20
##                      Asian
##                      95
```

```
##                                Black
##                                228
##                                Hispanic
##                                500
##                                More than one race
##                                81
## Native Hawaiian/Other Pacific Islander
##                                20
##                                White
##                                1470
table(pisa_train$grade) # ordered factor w at least 3 levels

##
##      8      9     10     11     12
##      2   188  1730   491      3
```

2.2) Unordered factors in regression models

To include unordered factors in a linear regression model, we define one level as the “reference level” and add a binary variable for each of the remaining levels. In this way, a factor with n levels is replaced by $n-1$ binary variables. The **reference level** is typically selected to be **the most frequently occurring level** in the dataset.

As an example, consider the unordered factor variable “color”, with levels “red”, “green”, and “blue”.

- If “green” were the reference level, then we would add binary variables “color_red” and “color_blue” to a linear regression problem.
 - All red examples would have color_red=1 and color_blue=0.
 - All blue examples would have color_red=0 and color_blue=1.
 - All green examples would have colorred=0 and colorblue=0.

Now, consider the variable “raceeth” in our problem, which has levels “American Indian/Alaska Native”, “Asian”, “Black”, “Hispanic”, “More than one race”, “Native Hawaiian/Other Pacific Islander”, and “White”. Because it is the most common in our population, we will select White as the reference level.

Which binary variables will be included in the regression model? (Since ‘white’ is the reference level, binary variables are every other race)

```
table(pisa_train$raceeth)

##
##      American Indian/Alaska Native
##                                20
##                                Asian
##                                95
##                                Black
##                                228
##                                Hispanic
##                                500
##                                More than one race
##                                81
## Native Hawaiian/Other Pacific Islander
##                                20
##                                White
##                                1470
```

2.3) Examle unordered factors

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.) (ans: Black, American Indian, Hispanic, More than one race, Native Hawaiian)

3.1) Building a model

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:

```
pisa_train$raceeth = relevel(pisa_train$raceeth, "White")
pisa_test$raceeth = relevel(pisa_test$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 = lm(Y ~ X1 + X2 + X3, data = Train)
```

You would use the following command to build your model:

```
LinReg = lm(Y ~ ., data = Train)
```

```
lm_score = lm(readingScore ~ ., data = pisa_train)
```

```
summary(lm_score)
```

```
##
## Call:
## lm(formula = readingScore ~ ., data = pisa_train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -247.44  -48.86    1.86   49.77  217.18
##
## Coefficients:
##              Estimate Std. Error
## (Intercept)   143.766333   33.841226
## grade          29.542707    2.937399
## male         -14.521653    3.155926
## raceethAmerican Indian/Alaska Native -67.277327   16.786935
## raceethAsian   -4.110325    9.220071
## raceethBlack  -67.012347    5.460883
## raceethHispanic -38.975486    5.177743
## raceethMore than one race -16.922522    8.496268
## raceethNative Hawaiian/Other Pacific Islander -5.101601   17.005696
## preschool     -4.463670    3.486055
## expectBachelors 55.267080    4.293893
```

```

## motherHS                6.058774    6.091423
## motherBachelors         12.638068    3.861457
## motherWork              -2.809101    3.521827
## fatherHS                4.018214    5.579269
## fatherBachelors         16.929755    3.995253
## fatherWork              5.842798    4.395978
## selfBornUS              -3.806278    7.323718
## motherBornUS            -8.798153    6.587621
## fatherBornUS            4.306994    6.263875
## englishAtHome           8.035685    6.859492
## computerForSchoolwork   22.500232    5.702562
## read30MinsADay         34.871924    3.408447
## minutesPerWeekEnglish   0.012788    0.010712
## studentsInEnglish       -0.286631    0.227819
## schoolHasLibrary        12.215085    9.264884
## publicSchool            -16.857475    6.725614
## urban                   -0.110132    3.962724
## schoolSize              0.006540    0.002197
##
## t value Pr(>|t|)
## (Intercept)             4.248 2.24e-05 ***
## grade                   10.057 < 2e-16 ***
## male                    -4.601 4.42e-06 ***
## raceethAmerican Indian/Alaska Native -4.008 6.32e-05 ***
## raceethAsian            -0.446 0.65578
## raceethBlack            -12.271 < 2e-16 ***
## raceethHispanic         -7.528 7.29e-14 ***
## raceethMore than one race -1.992 0.04651 *
## raceethNative Hawaiian/Other Pacific Islander -0.300 0.76421
## preschool              -1.280 0.20052
## expectBachelors         12.871 < 2e-16 ***
## motherHS                0.995 0.32001
## motherBachelors         3.273 0.00108 **
## motherWork              -0.798 0.42517
## fatherHS                0.720 0.47147
## fatherBachelors         4.237 2.35e-05 ***
## fatherWork              1.329 0.18393
## selfBornUS              -0.520 0.60331
## motherBornUS            -1.336 0.18182
## fatherBornUS            0.688 0.49178
## englishAtHome           1.171 0.24153
## computerForSchoolwork   3.946 8.19e-05 ***
## read30MinsADay         10.231 < 2e-16 ***
## minutesPerWeekEnglish   1.194 0.23264
## studentsInEnglish       -1.258 0.20846
## schoolHasLibrary        1.318 0.18749
## publicSchool            -2.506 0.01226 *
## urban                   -0.028 0.97783
## schoolSize              2.977 0.00294 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 73.81 on 2385 degrees of freedom
## Multiple R-squared:  0.3251, Adjusted R-squared:  0.3172
## F-statistic: 41.04 on 28 and 2385 DF, p-value: < 2.2e-16

```

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).

3.2) Computing the root-mean squared error of the model

```
SSE_train = sum(lm_score$residuals ^2)
RMSE_train = sqrt(SSE_train/ nrow(pisa_train))
```

```
SSE_train
```

```
## [1] 12993365
```

```
RMSE_train
```

```
## [1] 73.36555
```

3.3) Comparing predictions for similar students

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?

Use the coefficient for grade (29.542707) then multiply it by 2 => ansL 59.09

3.4) Interpreting model coefficients

What is the meaning of the coefficient associated with variable raceethAsian? Coeff = -4.110325 => **Predicted difference in the reading score between an Asian student and a white student who is otherwise identical

3.5) Identifying variables lacking statistical significance

Based on the significance codes, which variables are candidates for removal from the model? Select all that apply. (We'll assume that the factor variable raceeth should only be removed if none of its levels are significant.)

4.1) Predicting on unseen data

```
pred_test = predict(lm_score, newdata = pisa_test)
summary(pred_test)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   353.2   482.0   524.0   516.7   555.7   637.7
```

```
range = 637.7 - 353.2
range
```

```
## [1] 284.5
```


4.2) Test set SSE and RMSE

What is the sum of squared errors (SSE) of lmScore on the testing set? What is the root-mean squared error (RMSE) of lmScore on the testing set?

```
SSE_test = sum((pred_test - pisa_test$readingScore)^2)
```

```
RMSE = sqrt(SSE_test/nrow(pisa_test))
```

```
SSE_test
```

```
## [1] 5762082
```

```
RMSE
```

```
## [1] 76.29079
```

4.3) Baseline prediction and test-set SSE

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.

The baseline model is defined as the mean of the values seen in the training set (baseline prediction: the average value of dependent variables)

```
mean(pisa_train$readingScore)
```

```
## [1] 517.9629
```

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).

```
SSE_baseline = sum((mean(pisa_train$readingScore) - pisa_test$readingScore)^2)
```

```
SSE_baseline
```

```
## [1] 7802354
```

4.4) Test-set R-squared

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

```
SSE_test = sum((pisa_test$readingScore - pred_test)^2)
```

```
SST_test = sum((pisa_test$readingScore - mean(pisa_train$readingScore))^2)
```

```
R_squared_test = 1 - (SSE_test / SST_test)
```

```
SSE_test
```

```
## [1] 5762082
```

```
SST_test
```

```
## [1] 7802354
```

```
R_squared_test
```

```
## [1] 0.2614944
```