

Correlation and regression using the Lahman database for baseball

Michael Lopez, Skidmore College

Overview

In today's lab, we are going to explore parts of the Lahman package, while also linking to the statistical concepts in a bivariate analysis, including correlation, regression, and R-squared.

First, load the libraries.

```
library(Lahman)
library(tidyverse)
```

There is so much to the Lahman database, and in this course, we'll only touch the tip of the iceberg. First, scroll to the data frame list on the last page of the tutorial. Under **Data sets**, there is a list of roughly 25 data sets that come with the Lahman package.

You can also get a list of the data frames by entering the following command.

```
LahmanData
```

1. Which data frames in the Lahman package has the largest number of observations? Which have the fewest?

We're going to start by using the **Teams** data.

```
data(Teams)
head(Teams)
tail(Teams)
```

The **Teams** data contains team-level information for every year of baseball season, from 1871 - 2014. That's a lot of years! To make things a bit easier, we are going to focus on the modern era of baseball, which is generally considered to be the period from 1970 onwards.

Let's look at teams in the modern era, using the **filter()** command.

```
Teams_1 <- filter(Teams, yearID >= 1970)
head(Teams_1)
```

We store the newer data set as **Teams_1**, and you can tell **Teams_1** is a smaller data set because it begins in 1970, and not 1871.

Next, let's create a few missing team-level variables that we are going to need for the lab, using the **mutate()** command.

```
Teams_1 <- mutate(Teams_1, X1B = H - X2B - X3B - HR,
                  TB = X1B + 2*X2B + 3*X3B + 4*HR,
                  RC = (H + BB)*TB/(AB + BB),
                  RC.new = (H + BB - CS)*(TB + (0.55*SB))/(AB+BB) )
```

This creates four new team-level variables, **X1B** (number of singles), **TB** (total bases), **RC** (runs created), and **RC.new** (a newer runs created formula). The first formula for runs created is the basic one, discussed previously in lecture, while the second one uses a tweak for stolen bases.

Notice that in the above code, the data name (**Teams_1**) remained unchanged. Alternatively, we could have created a new name (say, **Teams_2**) if we had wanted.

Bivariate analysis

Correlation and scatter plots

Using `Teams_1`, let's quantify the strength, shape, and direction of the association between runs created and runs.

```
ggplot(data = Teams_1, aes(RC, R)) +  
  geom_point()  
  
Teams_1 %>%  
  summarise(cor_var = cor(R, RC))
```

2. Repeat the analysis above (which uses runs created) using the newer formula for runs created (`RC.new`)
3. Which variable - `RC` or `RC.new` - shows a stronger link with runs? What does this entail about the updated stolen bases formula?

Let's identify some additional ways of quantifying the association between two variables.

As an example, let's look at the relationship between a team's home runs and its runs.

```
Teams_1 %>%  
  summarise(cor_var = cor(R, HR),  
            r_sq = cor_var^2)
```

4. Interpret the correlation coefficient and the R-squared values between home runs and runs.

Simple linear regression.

From introductory statistics, you'll remember the simple linear regression (SLR) equation. Recall, here's the estimated SLR fit:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 * x_i$$

In the model above, \hat{y}_i represents our predicted value of an outcome variable given x_i , while $\hat{\beta}_0$ and $\hat{\beta}_1$ are the estimated intercepts and slopes, respectively.

In our example, we can plug in our variable names as follows:

$$\hat{R}_i = \hat{\beta}_0 + \hat{\beta}_1 * HR_i$$

5. Make a scatter plot of `R` as a function of `HR`. Using the function, make educated guesses for $\hat{\beta}_0$ and $\hat{\beta}_1$.

Fortunately, we don't have to make educated guesses. Let's look at the fit a model of runs as a function of home runs using the R command `lm()`.

```
fit_1 <- lm(R ~ HR, data = Teams_1)  
summary(fit_1)
```

The `summary()` code gives the regression output, as well as other model characteristics.

6. Write the estimated regression line. Does the actual estimated regression line resemble your guesses from question 5)?
7. Interpret both the intercept and the slope (for `HR`) in the estimated regression equation. Is the intercept a useful term here?
8. Estimate the number of runs that a team with 250 home runs would score. Alternatively, estimate how many home runs a team hit if they scored 575 runs.
9. Is the link between home runs and runs significant? How can you tell?

10. Use a similar code to determine if there is there a significant linear association between team runs (R) and the number of times that team was caught stealing (CS)?

Analyzing several variables simultaneously

We close the lab by looking at how one would analyze a subset of variables to judge which are most strongly associated with team success.

First, we use the `select()` command to reduce our data set from several dozen columns to only the ones we are interested in studying.

```
Teams_2 <- Teams_1 %>%  
  select(R, RC.new, RC, RA, HR, SO, attendance)
```

Next, we calculate the pairwise correlation coefficient between each variable.

```
cor_matrix <- cor(Teams_2)  
cor_matrix  
round(cor_matrix, 3)
```

11. Of the variables listed, which boast the strongest and weakest correlations with runs scored? Sidenote: What does the `round()` command do?

This is a good time to point out that, as is usually the case with observational data like this, strong links between two variables do not entail that one variable causes the other. While hitting home-runs will likely cause teams to score more runs, it is not necessarily the case that higher attendance causes teams to score more runs.

12. Think of one reason why attendance and runs are significantly correlated besides saying that attendance causes teams to score more runs.

Plotting correlations

There are lots of fun plots to make in R. Here are a few.

```
library(corrplot)  
corrplot(cor_matrix, method="number")  
corrplot(cor_matrix, method="circle", type = "lower")
```

Repeatability

In addition to wanting a metric to correlate with success and to reflect individual talent, it is worth looking at how well a metric can predict a future, unknown performance.

This requires some careful coding. Using the `dplyr` package, we create a new data frame, `Teams_3`, which arranges the team-level data by franchise and year, and then calculates the number of runs scored in the following year as the variable `next.R`.

```
Teams_3 <- Teams_1 %>%  
  arrange(franchID, yearID) %>%  
  group_by(franchID) %>%  
  mutate(next_R = lead(R))  
head(Teams_3)
```

13. Verify that the last column of `Teams.2` contains the future runs scored for each team in each row (hint: use the `select` command)

To close, we look at which team-level variables most strongly correlate with future runs scored.

```
Teams_3 <- Teams_3 %>% ungroup() ### Needed to remove grouping from earlier code
Teams_4 <- Teams_3 %>%
  select(next_R, R, RC, X1B, X2B, X3B, HR, SLG, OBP, OPS, attendance)
cor_matrix <- cor(Teams_4, use="pairwise.complete.obs")
cor_matrix
```

14. Which variables appear to be the best at predicting a team's runs scored in the following year? Which appears to be the worst?
15. Related: Compare the correlation of number of singles (X1B) to `next.R` and `R`, as well as number of home runs (HR) to `next.R` and `R`. Think carefully about how this would impact your recommendations of how to build a team. Which measures should be looked at most closely? Which measures appear to be mostly noise?
16. For fun: Go to the Shiny app here. Take a few guesses, and then click **View scatterplot of correlation between your guesses and actual correlation**. How good are you at guessing the correlation coefficient?