

Pitcher prediction using the Lahman database for baseball

Michael Lopez, Skidmore College

Overview

In this lab, we'll learn model comparison tools using multivariate regression. Next, we'll apply our tools to derive predictions of pitcher performance.

First, recall that we have to load the requisite libraries that we'll need (and we may have to install them, too. As always, once a package is downloaded, you do not need to run the `install.packages()` code again.

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

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

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

```
Teams.1 <- filter(Teams, yearID >= 1970)
head(Teams.1)
```

Comparing multiple regression models.

There's an old saying in statistics, attributed to George Box: *all models are wrong, some are useful*. In practice, we never know if our regression model is correctly specified; e.g., that it is really the case the y , x_1 , ... and x_{p-1} are linearly related. All we can do is hope...and try a few analytical tools.

Let's try to come up with a few models of RA: runs against. First, we start with a recap of multiple regression.

```
library(broom)
fit.1 <- lm(RA ~ HRA + BBA + SOA + HA + attendance, data = Teams.1)
tidy(fit.1)
```

1. Write the estimated model above.
2. Using the model in question (1), interpret the coefficient on HRA.
3. Using the model in (1), interpret the coefficient on **attendance**. Then come up with a better way to interpret the coefficient on **attendance**.
4. Remove HRA from the model and re-fit. Do your other coefficients change? Can you explain the difference?

After fitting a linear regression model, it is appropriate to check assumptions. First, we check the appropriateness of the normal distribution for residuals.

```
qqnorm(fit.1$resid)
qqline(fit.1$residuals)
```

Next, we compare the residuals to the fitted values, checking for the assumptions of independence among the residuals, as well as the constant variance assumption.

```
ggplot(data = fit.1, aes(x = fit.1$fitted.values, y = fit.1$residuals)) +
  geom_point() +
  geom_smooth()
```

5. What do the residual plots suggest about the assumptions of our linear regression model? What about the model makes it possibly a poor fit?

Speaking of residuals, let's take a deeper look at individual predictions.

The 1970 Atlanta Braves allowed 185 home runs, 478 walks, struck out 960 batters, and gave up 1451 hits. Their attendance was 1078848. As it turns out, the Braves are the first row of our data set.

```
Teams.1 %>%
  slice(1)

Teams.1 <- Teams.1 %>%
  mutate(fitted.value = predict(fit.1))

Teams.1 %>%
  slice(1)
```

6. The `predict` command calculates the fitted number of runs using a model such as ours (our model is `fit.1`). Using the code above, how many runs did our model predict that Braves to have allowed? What is the residual for the number of runs allowed by the Braves? Did our model overestimate or underestimate Atlanta's performance?

R-squared

There are lots of ways to measure the success of a regression model. The most common metric is **R-squared**, which you'll recall is the fraction of variability in the outcome which is explained by the regression model. Larger R-squared's are, in principal, better.

You can access the R-squared via the `summary` command:

```
summary(fit.1)
```

Using our model, we would interpret the R-squared as follows:

90.09% of the variability in the number of runs allowed by a team can be explained by the linear model with HRA, BBA, SOA, HA, and attendance.

While popular, the traditional R-squared is also flawed. Let's see how. In the following code, we'll create two new variables in the `Teams.1` data set, `rand1` and `rand2`, which are random normal variables.

```
set.seed(0)
Teams.1 <- mutate(Teams.1,
  rand1 = rnorm(nrow(Teams.1)),
  rand2 = rnorm(nrow(Teams.1)))
head(Teams.1)
```

Let's see what happens when we include `rand1` and `rand2` to our regression fit.

```
fit.2 <- lm(RA ~ HRA + BBA + SOA + HA + attendance + rand1 + rand2, data = Teams.1)
summary(fit.2)
```

Even when we added random noise to the model, R-squared went up!

That's not a good thing, at least when it comes to making model comparisons. In fact, it's a property of R-squared that, no matter what variable you add to a given model, the R-squared cannot go down. As a

result, R-squared is not useful for model comparisons, but more to gain a sense of how much of a drop in the variability in the outcome can be explained by the model's fit.

In place of R-squared, R also shows a formula for an adjusted R-squared, which penalizes models for adding unneeded parameters. However, this metric also has weaknesses.

7. When *would* be an appropriate time to compare R-squared's from two different models?
8. What other approaches to picking a model may be more appropriate than R-squared?