Temple University

SEI x Phillies 2024 Hackathon:

Estimating Strike Probabilities Through Umpire Tendencies

**CODE INSTRUCTIONS**

Kale Wiley, Gibson Hurst, Michael Mucciolo, & Cara Fredericks

April 24, 2024

Instructions for running code:

1. Import necessary libraries:
   a. library(readxl)
   b. library(dplyr)
   c. library(tidyr)
   d. library(ggplot2)
2. Import Dataset (We kept the name the same):
   a. baseball_copy <- read_excel("baseball copy.xlsx")
3. Run regression code (highlighted chunk):

```
104
105  # Running regression on each variable against description
106  # (seeing how much of an impact each variable has on the call
107  # the umpire calls for each pitch)
108
109  # Converting description column variables (1 for called_strike and 0 for ball)
110  baseball_copy$description_binary <- ifelse(baseball_copy$description == "called_strike", 1, 0)
111
112  # Regression for each variable against description
113  variables <- c("pitch_type", "stand", "p_throws", "balls", "strikes", "plate_x", "plate_z", "sz_top", "sz_bot", "zone")
114  results <- lapply(variables, function(var) {
115    formula <- as.formula(paste("description_binary ~", var))
116    first_model <- glm(formula, data = baseball_copy, family = binomial)
117    return(summary(first_model)$coefficients)
118  })
119
120  names(results) <- variables
121  print(results)
122
123
124  # Creating second logistic regression (simpler to use for manipulation, same principles as first regression)
125  logmodel <- glm(description_binary ~ ., data = baseball_copy[, c("description_binary", variables)], family = binomial)
126
127  summary(logmodel)
```

4. Run accuracy analysis:

```
# Accuracy analysis of regression 2 to determine which cutoff level should be used
# to make model as accurate as possible (for prediction analysis)

# Defining cutoff levels to determine the threshold probability for classifying an outcome as either positive (strike) or negative (ball)
cutoffs <- seq(0.1, 0.9, 0.1)

accuracy <- NULL

# Calculating accuracy for each cutoff level and creating plot
for (i in seq_along(cutoffs)){
  prediction <- ifelse(logmodel$fitted.values >= cutoffs[i], 1, 0)
  accuracy <- c(accuracy, sum(prediction == baseball_copy$description_binary)/length(prediction)*100)
}

plot(cutoffs, accuracy, type='l', xlab="Cutoff Level", ylab="Accuracy %",
     main="Cutoff Level vs. Model Accuracy")
```

5. Testing model (adjust variables however you please):

```r
152
153
154    # # Testing the model (input variables below in new_pitch data frame and
155    # input cutoff level in ifelse prediction section, output will be the prediction
156    # of whether the pitch would be a ball or a strike given all variables)
157
158    # Running new_pitch data frame to analyze accuracy of regression model
159    # (based on variables input should the pitch be called a ball or a strike)
160    new_pitch <- data.frame(
161      pitch_type = "SL",
162      stand = "L",
163      p_throws = "L",
164      balls = 3,
165      strikes = 2,
166      plate_x = 0.5,
167      plate_z = 1.5,
168      sz_top = 2.5,
169      sz_bot = 1.0,
170      zone = 4
171    )
172
173    # Prediction using the regression model
174    prediction <- predict(logmodel, newdata = new_pitch, type = "response")
175
176    # Determining if the pitch will be called a strike or ball based on the cutoff level and other variables
177    ifelse(prediction >= 0.8, "Strike", "Ball")
178    |
179
178:1    (Top Level)                                                    R S
```

Console