

## Derin Gezgin | Camel ID:00468038

**Problem:** For the MLBStandings2016 data, do the following:

- a. [35 points] Regress WinPct on ERA and League and report the fitted model. Interpret the regression coefficient for the League predictor. Make a plot of ERA vs WinPct with separate lines for the two leagues. Is League a significant predictor of WinPct in the presence of ERA. Show hypothesis test for the League predictor.

*Regress WinPct on ERA and League and report the fitted model.*

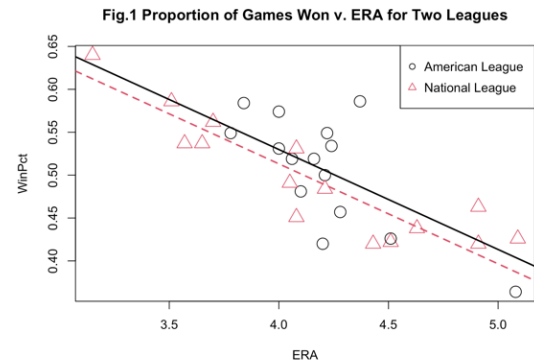
*Make a plot of ERA v. WinPct with separate lines per category.*

$$\hat{Y} = 0.996 - 0.117x_1 - 0.017x_2 \text{ where:}$$

$Y \rightarrow$  Proportion of games won (WinPct)

$x_1 \rightarrow$  Earned run average (earned runs allowed per 9 innings)

$x_2 \rightarrow$  League: AL=American or NL=National



*Interpret the regression coefficient for the League predictor*

In the **R** analysis, the **American League** is the base category with an  $x_2$  value of 0.

In this case, the fitted model for the American League ( $x_2 = 0$ ) would be;

$$\hat{Y} = 0.996 - 0.117x_1$$

**0.996** is the estimated proportion of games won, for the American League when the earned run average is 0.

*On the other hand, the **National League** is the other category with an  $x_2$  value of 1.*

In this case, the fitted model for the National League ( $x_2 = 1$ ) would be;

$$\hat{Y} = 0.979 - 0.117x_1$$

**0.979** is the estimated proportion of games won, for the National League when the earned run average is 0.

*As the intercept of the American League is higher, we can say that -in terms of proportion of games won- the American League performs better than the National League*

### **General Interpretation of the Regression Coefficients**

$\beta_1 = 0.117 \rightarrow$  Estimated change in proportion of games won for a unit increase in earned run average for either league.

$\beta_2 = 0.017 \rightarrow$  Estimated difference in proportion of games won for the American league, as compared to the national league for any earned run average.

$\beta_0 = 0.996 \rightarrow$  The estimated proportion of games won for the American League when the earned run average is 0.

$\beta_0 + \beta_2 = 0.979 \rightarrow$  The estimated proportion of games won for the National League when the earned run average is 0.

***Is League a significant predictor of WinPct in the presence of ERA. Show hypothesis test for the League predictor.***

***Step 1: Formulate Null and Alternative Hypothesis***

$H_0: \beta_2 = 0 \rightarrow$  League is not a significant predictor of proportion of games won (WinPct).

$H_A: \beta_2 \neq 0 \rightarrow$  League is a significant predictor of proportion of games won (WinPct).

***Step 2: Set level of significance:*** 0.05

***Step 3: Test statistic:*** -1.135

***Step 4: P-Value:*** 0.266

***Step 5: Conclusion:***

***P-Value*** is larger than our level of significance. This means that at 5% level of significance, we cannot reject the null hypothesis. In other words, *league* is not a significant predictor of proportion of games won (WinPct) in the presence of ERA.

- b. [25 points] For the model regressing WinPct on the predictors ERA and League include the interaction term. Report the fitted model and interpret the coefficients. Make a plot for the two leagues with different intercepts and slopes. Check if the interaction term is statistically significant.**

***Report the fitted model***

$$\hat{Y} = 1.149 - 0.153x_1 - 0.213x_2 + 0.047x_1x_2$$

$Y \rightarrow$  Proportion of games won (WinPct)

$x_1 \rightarrow$  Earned run average (earned runs allowed per 9 innings)

$x_2 \rightarrow$  League: AL=American or NL=National

$x_1x_2 \rightarrow$  Interaction term between ERA and League.

***Fitted Model for the American League:***  $x_2$  value of 0

$$\hat{Y} = 1.149 - 0.153x_1$$

***Fitted Model for the National League:***  $x_2$  value of 1

$$\hat{Y} = 0.936 - 0.106x_1$$

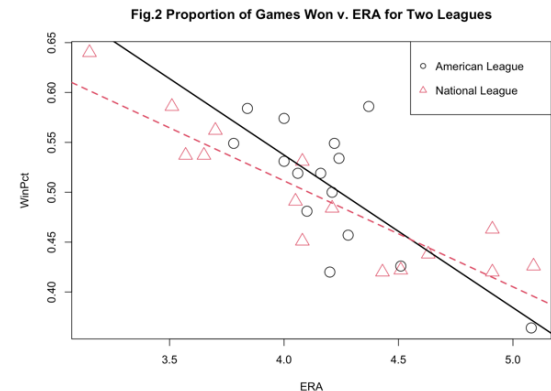
***Interpretation of the Regression Coefficients***

$\beta_1 = 0.153 \rightarrow$  Estimated change in proportion of games won for a unit increase in earned run average for the American league.

$\beta_2 = 0.213 \rightarrow$  Estimated difference in proportion of games won for the American league, as compared to the national league for a run average of 0.

$\beta_3 = 0.047 \rightarrow$  Estimated difference of proportion of games won for an increase in earned run average, for the American league compared to the National league.

$\beta_1 + \beta_3 = 0.106 \rightarrow$  Estimated change in proportion of games won for a unit increase in earned run average for the National League.



$\beta_0 = 1.149 \rightarrow$  The estimated proportion of games won for the American League when the earned run average is 0.

$\beta_0 + \beta_2 = 0.936 \rightarrow$  The estimated proportion of games won for the National League when the earned run average is 0.

### ***Checking if the Interaction Term is Statistically Significant***

#### ***Step 1: Formulate Null and Alternative Hypothesis***

$H_0: \beta_3 = 0 \rightarrow$  Interaction Term is not a significant predictor.

$H_A: \beta_3 \neq 0 \rightarrow$  Interaction term is a significant predictor

***Step 2: Set level of significance:*** 0.05

***Step 3: Test statistic:*** 1.185

***Step 4: P-Value:*** 0.247

***Step 5: Conclusion:***

*P-Value* is larger than our level of significance. This means that at 5% level of significance, we do not have enough evidence to reject the null hypothesis and conclude that the interaction term between *League* and earned run average is not a significant predictor of proportion of games won (WinPct).

- c. [25 points] Regress WinPct on the predictors ERA and Runs including their interaction term. Report the fitted model and interpret the coefficients. Check if the interaction term is statistically significant.

#### ***Report the Fitted Model***

$$\hat{Y} = 0.464 - 0.073x_1 + 0.0006x_2 - 0.00004x_1x_2$$

$$\hat{Y} = 0.464 - (0.073 - 0.00004x_2)x_1 + 0.0006x_2$$

$Y \rightarrow$  Proportion of games won

$x_1 \rightarrow$  Earned run average (earned runs allowed per 9 innings)

$x_2 \rightarrow$  Runs (Number of runs scored)

$x_1x_2 \rightarrow$  Interaction term between earned run average and runs scored.

#### ***Interpretation of the Coefficients***

$\beta_0 = 0.464 \rightarrow$  Estimated proportion of games won when the earned run average and number of runs scored are both 0.

$\beta_1 = -0.073 \rightarrow$  Estimated change in proportion of games won for a unit increase in earned run average when number of runs scored is 0.

$\beta_2 = 0.0006 \rightarrow$  Estimated change in proportion of games won for a unit increase in number of runs scored when earned run average is 0.

$\beta_3 = -0.00004 \rightarrow$  Estimate of the modification to the change in proportion of games won for a unit increase in earned run average in case of a certain number of runs scored.

***Checking the statistical significance of the interaction term***

***Step 1: Formulate Null and Alternative Hypothesis***

$H_0: \beta_3 = 0 \rightarrow$  Interaction Term is not a statistically significant predictor.

$H_A: \beta_3 \neq 0 \rightarrow$  Interaction term is a statistically significant predictor

***Step 2: Set level of significance ( $\alpha$ ): 0.05***

***Step 3: Test statistic: -0.198***

***Step 4: P-Value: 0.845***

***Step 5: Conclusion:***

*P-Value* is larger than our level of significance. At 5% level of significance, we do not have enough evidence to reject the null hypothesis. In other words, this means that at 5% level of significance, the interaction term between ERA and Runs is not a significant predictor of WinPct.

- d. [15 points] Perform ANOVA with FM as model with ERA, Runs, and their interaction term and RM as the ERA and Runs model.

***Step 1: Formulate Null and Alternative Hypothesis***

$H_0: \beta_3 = 0 \rightarrow$  Interaction Term is not a statistically significant predictor. Favoring the removed model

$H_A: \beta_3 \neq 0 \rightarrow$  Interaction term is a statistically significant predictor. Favoring the full model.

***Step 2: Set level of significance ( $\alpha$ ): 0.05***

***Step 3: Test statistic: 0.039***

***Step 4: P-Value: 0.8449***

***Step 5: Conclusion:***

*P-Value* is larger than our level of significance. This means that at 5% level of significance, we do not have evidence to reject the null hypothesis. We can conclude the interaction term is not a statistically significant predictor and we would favor the removed model which does not include the interaction term.