

# Predicting Baseball Pitcher's Earned Run Average (ERA)

Aaron Wu, Ruchi Tiwari, Edwin Muñoz, Shakty Juarez

December 2024

## Problem Statement

The primary objective of this project is to develop a predictive regression model that accurately forecasts a baseball pitcher's Earned Run Average (ERA) based on various performance metrics. With a dataset encompassing 106 potential predictors, the challenge lies in identifying the most significant regressors that contribute meaningfully to the prediction of ERA. This study aims to:

- Enhance the understanding of factors influencing pitcher performance.
- Provide a robust analytical tool for evaluating and forecasting ERA.
- Address multicollinearity and model overfitting using advanced regression techniques.

## Introduction

Baseball analytics increasingly leverages statistical methodologies to evaluate player performance and inform strategic decisions. Among key performance indicators, ERA serves as a critical measure of a pitcher's effectiveness, representing the average number of earned runs allowed per nine innings pitched. Accurate prediction of ERA is invaluable for team management, player development, and fantasy sports enthusiasts.

This project aims to construct a regression model that predicts a pitcher's ERA using a comprehensive set of performance metrics. The dataset includes variables such as:

- Strike percentage (`k_percent`)
- Walk percentage (`bb_percent`)
- Whiff percentage (`whiff_percent`)
- Weighted on-base average (`woba`)
- Expected batting average (`xba`), among others.

The study's goals were:

1. Identify the most significant predictors of ERA.
2. Develop a model balancing accuracy and interpretability.
3. Address multicollinearity and overfitting with advanced techniques.
4. Validate the model using diagnostic tests for reliability.

### **Key Questions:**

1. Which performance metrics are the most significant predictors of ERA?
2. How do advanced regression techniques (Ridge, LASSO, Elastic Net) improve model performance?
3. How can multicollinearity issues be mitigated?
4. How well does the final model generalize to new data?

## **Data Collection**

The analysis utilizes two primary datasets:

- **Pitching Data 2023** (`pitching_data_2023.csv`): Comprehensive pitching statistics for the 2023 season.

- **Additional Stats** (`additional_stats.csv`): Supplementary statistics for deeper insights.

Datasets were merged using the unique identifier `last_name.first_name`. Preliminary data exploration ensured proper alignment and integrity. Summary statistics and visualizations were employed to understand variable distributions, highlight outliers, and identify potential multicollinearity issues.

## Data Analysis and Results

Multiple linear regression models tested specific hypotheses:

**Model 1:**  $p\_era \sim k\_percent + bb\_percent$

*Hypothesis: Basic stats like strikeout and walk rates affect performance.*

**Model 2:**  $p\_era \sim woba + whiff\_percent$

*Hypothesis: Preventing strong batting and inducing misses impacts effectiveness.*

**Model 3:**  $p\_era \sim player\_age + xba \times woba$

*Hypothesis: Age and expected batting performance influence ERA.*

**Model 4:**  $p\_era \sim swing\_percent + hard\_hit\_percent + whiff\_percent$

*Hypothesis: Swing behavior and contact quality explain ERA.*

**Model 5:**  $p\_era \sim player\_age \times k\_percent + bb\_percent$

*Hypothesis: Control, strikeouts, and age impact ERA.*

### Model Evaluation:

- Residual vs. fitted plots indicated no clear patterns.
- AIC values identified **Model 2** (`woba + whiff_percent`) as the best fit with the lowest AIC (91.89).
- Coefficient plots revealed `woba` as a major predictor, with minimal impact from `player_age` and `swing_percent`.

# Regularization Techniques

To address multicollinearity:

- **Ridge Regression:** Stabilized coefficients and optimized  $\lambda$  via cross-validation.
- **LASSO Regression:** Selected key predictors by shrinking weaker coefficients to zero.
- **Elastic Net:** Combined Ridge and LASSO penalties for balance and reduced overfitting.

# Diagnostics and Validation

1. **Residual Analysis:** Plots showed no significant patterns.
2. **QQ-Plot:** Confirmed normality of residuals.
3. **Influence Measures:** Identified influential outliers (e.g., Shohei Ohtani, Blake Snell).
4. **Autocorrelation:** Durbin-Watson test confirmed no autocorrelation.
5. **VIF Analysis:** Reduced multicollinearity post-regularization.

# Summary and Discussion

This project successfully developed a robust regression model for predicting ERA, addressing multicollinearity and overfitting using advanced techniques. The final Elastic Net model achieved strong predictive capabilities, identifying `whiff_percent` and `woba:xba` interaction as key predictors.

# Future Work

- Expand datasets to include multiple seasons.
- Incorporate additional metrics (e.g., pitch velocity, movement).

- Refine parameter tuning via granular grid search.
- Address outliers through specialized modeling approaches.

## Appendix

```

1 # This part of the code helps to determine which regressors out of
  the hundred in
2 # the dataset showed some level of promise in creating a regression
  model that could
3 # accurately predict pitcher ERA
4 library(tidyverse)
5 library(stats)
6
7 file_path <- "/Users/aaronwu/Desktop/stat410 copy/final/pitching_
  data_2023.csv"
8 file_path1 <- "/Users/aaronwu/Desktop/stat410 copy/final/additional_
  stats.csv"
9 data <- read.csv(file_path)
10 data1 <- read.csv(file_path1)
11 merged_data <- merge(data, data1, by = "last_name..first_name")
12
13 colnames(merged_data)
14 head(merged_data)
15
16 models <- list(
17   lm_1 = lm(p_era ~ k_percent + bb_percent, data = merged_data),
18   lm_2 = lm(p_era ~ woba + whiff_percent, data = merged_data),
19   lm_3 = lm(p_era ~ player_age + xba * woba, data = merged_data),
20   lm_4 = lm(p_era ~ swing_percent + hard_hit_percent + whiff_percent
    , data = merged_data),
21   lm_5 = lm(p_era ~ player_age * k_percent + bb_percent, data =

```

```

    merged_data)
22 )
23
24 model_summaries <- lapply(models, summary)
25 model_aics <- sapply(models, AIC)
26 print(model_summaries)
27 print(model_aics)
28
29 par(mfrow = c(2, 2))
30 plot(models$lm_1)
31
32 coef_plots <- lapply(models, function(model) {
33   coef_df <- broom::tidy(model)
34   ggplot(coef_df, aes(x = term, y = estimate, fill = estimate > 0))
35     +
36     geom_col() +
37     coord_flip() +
38     labs(title = paste("Coefficients of", deparse(model$call[[2]])),
39           x = "Terms",
40           y = "Estimates")
41 })
42
43 resid_fitted_plots <- lapply(models, function(model) {
44   ggplot(model, aes(.fitted, .resid)) +
45     geom_point() +
46     geom_smooth(method = "lm", col = "red") +
47     labs(title = paste("Residuals vs. Fitted for", deparse(model$
48       call[[2]])),
49           x = "Fitted values",
50           y = "Residuals")
51 })

```

```
51 print(resid_fitted_plots)
```

Listing 1: Initial Data Analysis and Model Comparison

```
1 # This is the first model. We are trying to generate a linear
  regression model.
2 library(ggplot2)
3 library(broom)
4 library(dplyr)
5
6 model <- lm(p_era ~ k_percent + bb_percent + whiff_percent + woba +
  xba:woba,
7           data = merged_data)
8 model_summary <- summary(model)
9 print(model_summary)
10
11 tidy_model <- tidy(model)
12 print(tidy_model)
13
14 ggplot(tidy_model, aes(x = term, y = estimate, fill = estimate > 0))
  +
15   geom_col() +
16   coord_flip() +
17   labs(title = "Coefficients of the Linear Regression Model",
18        x = "Terms", y = "Estimates") +
19   scale_fill_manual(values = c("red", "blue"),
20                     labels = c("Negative", "Positive"))
```

Listing 2: Initial Linear Regression Model

```
1 # Look at the VIF to determine if methods like Ridge would be useful
2 library(car)
3
4 model <- lm(p_era ~ k_percent + bb_percent + whiff_percent + woba +
```

```

    xba * woba,
5         data = merged_data)
6 vif_values <- vif(model)
7 print(vif_values)
8 # In this case, high VIF for woba, xba, and woba:xba, so we try

```

Listing 3: VIF Analysis

```

1 # This is a ridge prediction. Is there any other way we can optimize
  this?
2 library(glmnet)
3
4 predictors <- model.matrix(p_era ~ k_percent + bb_percent + whiff_
  percent +
5                           woba + xba * woba, data = merged_data)[,-1]
6 response <- merged_data$p_era
7 lambda_values <- 10^seq(10, -2, length = 100)
8
9 ridge_model <- glmnet(predictors, response, alpha = 0,
10                      lambda = lambda_values, standardize = TRUE)
11 cv_ridge <- cv.glmnet(predictors, response, alpha = 0,
12                      type.measure = "mse", nfolds = 10)
13 plot(cv_ridge)
14
15 best_lambda <- cv_ridge$lambda.min
16 best_ridge_model <- glmnet(predictors, response, alpha = 0,
17                          lambda = best_lambda, standardize = TRUE)
18 coef(best_ridge_model)
19 # From the results, it seems fine, but we need to check for
20 # residuals/other techniques to validate model

```

Listing 4: Ridge Regression Analysis

```

1 predicted_values <- predict(best_ridge_model, s = "lambda.min", newx

```



```

    = predictors)
2 residuals <- response - predicted_values
3
4 plot(residuals, type = 'p', main = "Residual Plot",
5       xlab = "Predicted Values", ylab = "Residuals")
6 abline(h = 0, col = "red")
7 #Residual plot looks good, let's look at other plots

```

Listing 5: Ridge Regression Residual Analysis

```

1 # Trying other techniques like looking at QQ-plot, leverage points,
2 # autocorrelation checks, Durbin-Watson and Ljung-Box tests
3 library(car)
4 library(lmtest)
5 library(stats)
6
7 qqnorm(residuals)
8 qqline(residuals, col = "red")
9
10 influencePlot(model, id.method="identify", main="Influence Plot",
11               sub="Circle size is proportional to Cook's Distance")
12
13 acf(residuals(model), main="ACF of Residuals")
14
15 dw_result <- dwtest(model)
16 print(dw_result)
17
18 lb_test <- Box.test(residuals(model), type = "Ljung-Box")
19 print(lb_test)
20
21 # Overall, all of the techniques showed that the model we created
22 # was good.
23 # Out of curiosity, I want to examine what other models look like.

```

```

23 # First, looking at making some of the coefficients polynomial,
24 # and second looking at a LASSO regression

```

Listing 6: Model Diagnostics

```

1 # Make some coefficients polynomial
2 library(stats)
3
4 poly_model <- lm(p_era ~ poly(k_percent, 2) + poly(bb_percent, 2) +
5                   poly(whiff_percent, 2) + poly(woba, 2) + poly(xba,
6                   2) +
7                   poly(woba, 1)*poly(xba, 1), data = merged_data)
8 summary(poly_model)
9 plot(poly_model)
10
11 # Polynomial slightly fits the data better, but at the cost of
12 # overfitting/homoscedacity. Overall tradeoff is that it can capture
13 # more
14 # complex relationships, but is less robust for predicting on newer
15 # data

```

Listing 7: Polynomial Regression Analysis

```

1 # LASSO regression
2 library(glmnet)
3
4 predictors <- model.matrix(p_era ~ k_percent + bb_percent + whiff_
5                             percent +
6                             woba + xba * woba, data = merged_data)
7
8 response <- merged_data$p_era
9
10 set.seed(123)
11 lasso_model <- glmnet(predictors, response, alpha = 1,

```

```

10             lambda = 10^seq(4, -2, length = 100))
11 cv_lasso <- cv.glmnet(predictors, response, alpha = 1)
12 plot(cv_lasso)
13
14 best_lambda <- cv_lasso$lambda.min
15 best_lasso_model <- glmnet(predictors, response, alpha = 1, lambda =
    best_lambda)
16 #print(coef(best_lasso_model))
17
18 # LASSO shows that there is actually two variables, whiff_percent
    and woba: xba,
19 # that are causing potentially a more complex model than necessary.
20 # Because of that, let's try combining the LASSO and Ridge into an
21 # Elastic Net Regression

```

Listing 8: LASSO Regression Analysis

```

1 # Elastic Net Regression
2 library(glmnet)
3
4 predictors <- model.matrix(p_era ~ k_percent + bb_percent + whiff_
    percent +
5             woba + xba * woba, data = merged_data)
    [, -1]
6 response <- merged_data$p_era
7
8 set.seed(123)
9 cv_model <- cv.glmnet(predictors, response, alpha = 0.5, family = "
    gaussian",
10             standardize = TRUE, type.measure = "mse",
    nfolds = 10)
11 plot(cv_model)
12

```

```

13 best_lambda <- cv_model$lambda.min
14 final_model <- glmnet(predictors, response, alpha = 0.5,
15                       lambda = best_lambda, standardize = TRUE)
16 print(coef(final_model))
17
18 # This model looks super good. However, we need to validate it one
   more time.

```

Listing 9: Elastic Net Regression Analysis

```

1 library(glmnet)
2 library(car)
3 library(caret)
4 library(Metrics)
5
6 predictions <- predict(final_model, s = best_lambda, newx =
   predictors)
7 r_squared <- cor(predictions, response)^2
8 residuals <- response - predictions
9
10 plot(predictions, residuals, main = "Residuals vs Fitted",
11       xlab = "Fitted Values", ylab = "Residuals")
12 abline(h = 0, col = "red")
13
14 qqnorm(residuals)
15 qqline(residuals, col = "red")
16
17 standard_model <- lm(p_era ~ .,
18                     data = as.data.frame(cbind(p_era = response,
   predictors)))
19 vif_values <- vif(standard_model)
20 print(vif_values)
21

```

```

22 plot(predictions, sqrt(abs(residuals)), main = "Scale-Location Plot"
    ,
23       xlab = "Fitted Values", ylab = "Sqrt(|Residuals|)")
24 abline(h = 0, col = "red")
25
26 leverage_values <- hatvalues(standard_model)
27 plot(leverage_values, main = "Leverage Plot")
28 abline(h = 2 * mean(leverage_values), col = "red", lty = 2)
29
30 cooks_d <- cooks.distance(standard_model)
31 plot(cooks_d, main = "Cook's Distance", type = "h")
32 threshold <- 4 / (length(response) - length(coef(standard_model)))
33 abline(h = threshold, col = "red", lty = 2)
34
35 cat("R-squared:", r_squared, "\n")
36
37 # Overall, this model is probably the best we will achieve. Only
    thing is
38 # potentially looking at Cook's Distance because there are 4
    outliers.

```

Listing 10: Final Model Validation

```

1 # I just want to check outliers
2 library(glmnet)
3 library(dplyr)
4
5 predictors <- model.matrix(p_era ~ k_percent + bb_percent + whiff_
    percent +
6                               woba + xba * woba, data = merged_data)
    [, -1]
7 response <- merged_data$p_era
8

```

```

9 set.seed(123)
10 cv_model <- cv.glmnet(predictors, response, alpha = 0.5, family = "
    gaussian")
11 best_lambda <- cv_model$lambda.min
12 final_model <- glmnet(predictors, response, alpha = 0.5, lambda =
    best_lambda)
13
14 cooks_values <- cooks.distance(lm(p_era ~ .,
15                                 data = as.data.frame(cbind(p_era =
    response,
16                                                         predictors
    ))))
17 leverage_values <- hatvalues(lm(p_era ~ .,
18                                data = as.data.frame(cbind(p_era =
    response,
19                                                         predictors))
    ))
20
21 cooks_threshold <- 4 / (nrow(predictors) - ncol(predictors) - 1)
22 high_cooks <- which(cooks_values > cooks_threshold)
23
24 leverage_threshold <- 2 * ncol(predictors) / nrow(predictors)
25 high_leverage <- which(leverage_values > leverage_threshold)
26
27 influential_points <- sort(unique(c(high_cooks, high_leverage)))
28 influential_data <- merged_data[influential_points, ]
29 print(influential_data)
30
31 plot(response, pch = 20,
32       col = ifelse(seq_along(response) %in% influential_points, "red"
    , "black"))
33 legend("topright", legend = c("Influential", "Not influential"),

```

```

34         col = c("red", "black"), pch = 20)
35
36 # Most of these players have a few things in common:
37 # 1. High k_percent and bb_percent like Trevor Williams, Michael
    Kopech, Luke Weaver
38 # 2. Extreme innings pitched like Shohei Ohtani, Blake Snell
39 # 3. Combination of these 2: Spencer Strider

```

Listing 11: Outlier Analysis

## Self Reflections

**Aaron Wu** - Throughout this project, I dedicated approximately 15 hours to data collection, analysis, and report preparation. This endeavor significantly enhanced my understanding of regression modeling, particularly in handling multicollinearity and implementing regularization techniques. One major challenge was managing the extensive number of predictors, which necessitated a methodical approach to variable selection and model optimization. I learned the importance of thorough data exploration and the role of diagnostic tests in validating model assumptions. Encountering influential outliers underscored the need for vigilance in model interpretation and the consideration of exceptional cases. If I were to undertake this project again, I would allocate more time to exploring interactive visualization tools to gain deeper insights into data patterns. Additionally, incorporating domain expertise earlier in the process could guide more informed variable selection and modeling decisions. My advice to future students is to embrace a structured approach to data analysis, remain open to iterative model refinement, and prioritize understanding the underlying assumptions of each statistical method employed. Effective time management and continuous learning are key to navigating complex analytical projects successfully.

**Ruchi Tiwari** - Over the course of 2–3 weeks, I dedicated approximately 15 hours to the project, with much of the initial time spent identifying a suitable dataset and establishing a clear direction. Through this process, I learned that teamwork and consistent communication are critical to ensuring everyone remains aligned, especially during complex phases

of the project. A strong grasp of domain knowledge also proved essential for interpreting model results, managing high-leverage points or outliers, and making informed decisions about which variables to retain, transform, or interact within the model. I also realized the importance of evaluating multicollinearity among predictors and ensuring assumptions like linearity and homoscedasticity were met to maintain the reliability of the model. While we conducted thorough validation and testing, future iterations could benefit from additional diagnostic checks, such as cross-validation and sensitivity analyses, to strengthen confidence in the model's robustness. My key advice for future students is to start early, as the final weeks can become particularly demanding, and to review slide decks ahead of time to deepen your understanding of validation and modeling techniques before they're covered in class. Taking these steps, along with dedicating time to thoroughly explore your data and its statistical properties, can make the project process more manageable and ultimately more rewarding.

**Edwin Muñoz** - I spent about 15 hours working throughout the course of the project. Most of that time was spent validating the models to see which fit the problem best according to different metrics. This project taught me how to choose a new model based on the previous model's performance, which metrics to look at, what trends to look for, and how to do this all in a team. Although we were pretty happy with our result, we had a bit of a struggle dealing with the outliers', more specifically how they made it difficult to measure performance and to chose the next model. Were we to do this again, it would be good to include further diagnostic checks to reassure how well our model did. My biggest advice to future students would be to start early and to take a good look at your data before making any call, as it is easy to get behind or start on the wrong path.

**Shakty Juarez** - I dedicated approximately 15 hours to this project, focusing on designing and implementing the regularization techniques, including Ridge, LASSO, and Elastic Net regression. This involved parameter tuning, selecting optimal lambda values through cross-validation, and ensuring that significant predictors were retained while mitigating multicollinearity. Through this group project, I gained a deeper understanding of how regularization techniques stabilize models and improve predictive accuracy. Working with baseball



data also deepened my appreciation for how statistical methods can be applied to real-world sports analytics specifically in baseball player performance. If I were to revisit this project, I would incorporate additional data sets to potentially explore more interaction terms or nonlinear transformations to capture more complex relationships in predicting player woba. My advice to future students is to dedicate time to researching your options for advanced model validation techniques thoroughly, and find ones that fit your specific dataset the best.