

Forecasting MLB Game Attendance

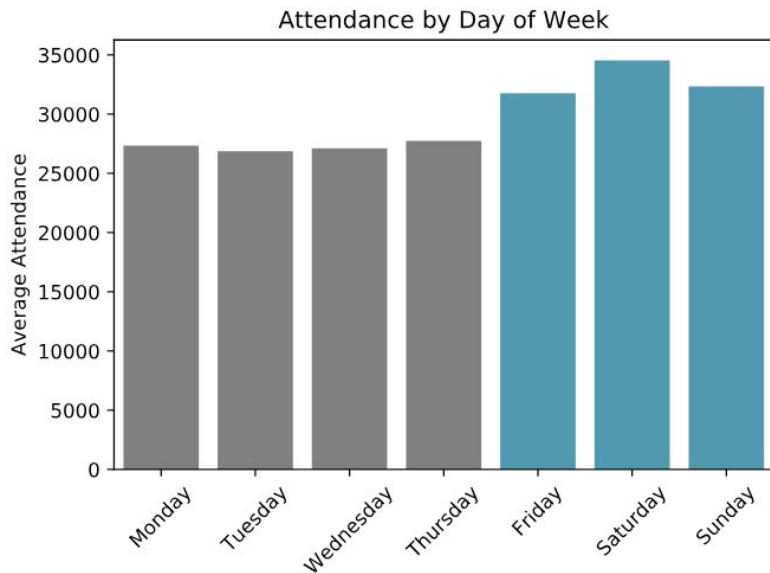
Rob Pagano

Predicting Attendance – Why it Matters

- Team marketing
- Local businesses
- Public Services

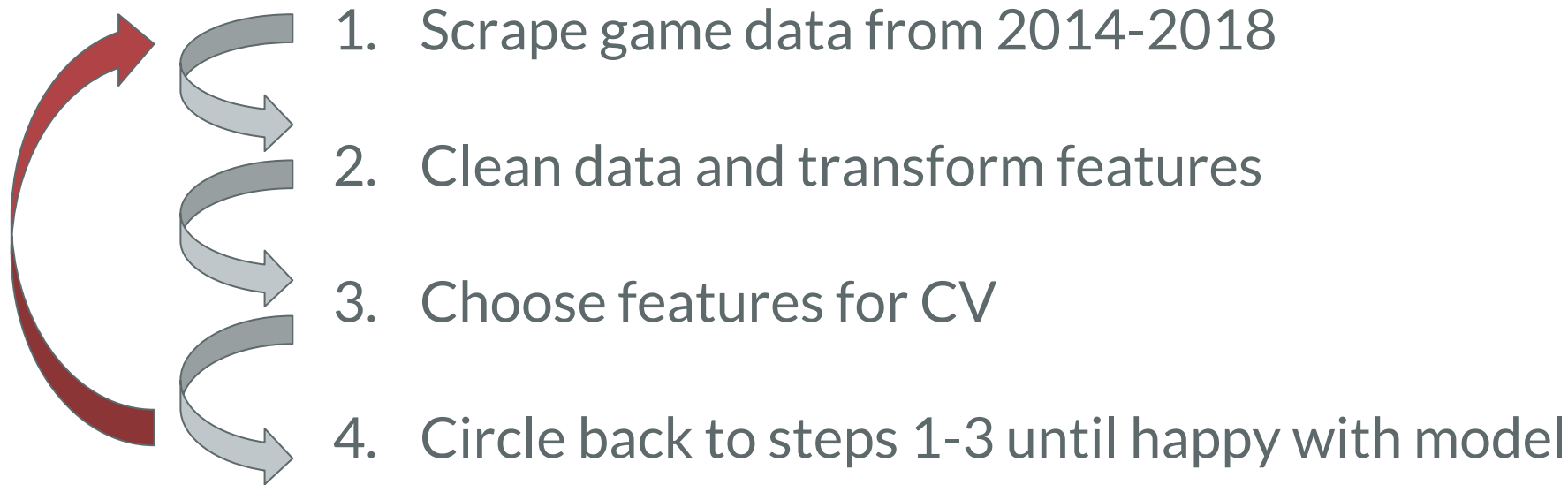


Approach



- What makes a *team* exciting
 - Baseball stats
 - Playoff contention
- What makes a *specific game* popular
 - Day of the week / Time of day
 - Opposing team
 - Promotions

Process



Feature Editing

Original OLS

- Rank
- Games back
- Night or day
- Wins last 10
- Mean runs last 10
- Run differential

RMSE \approx 9070.96

$R^2 \approx 0.16$

Feature Editing

Original OLS

- Rank
- Games back
- Night or day
- Wins last 10
- Mean runs last 10
- Run differential

RMSE \approx 9070.96

$R^2 \approx 0.16$



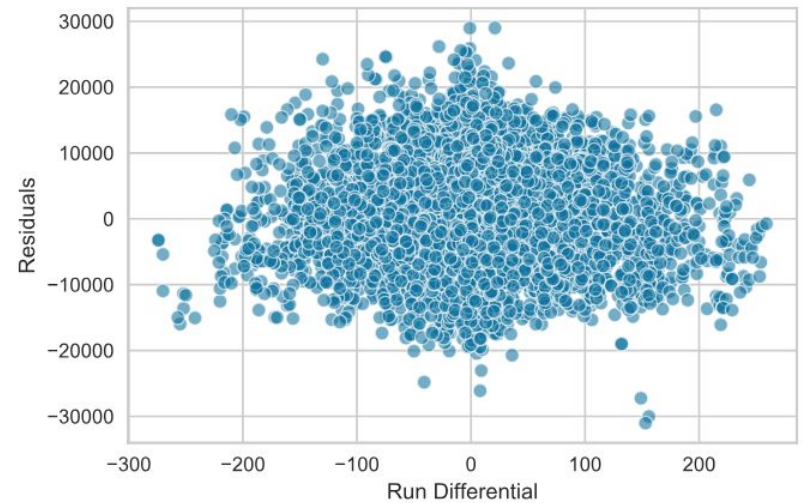
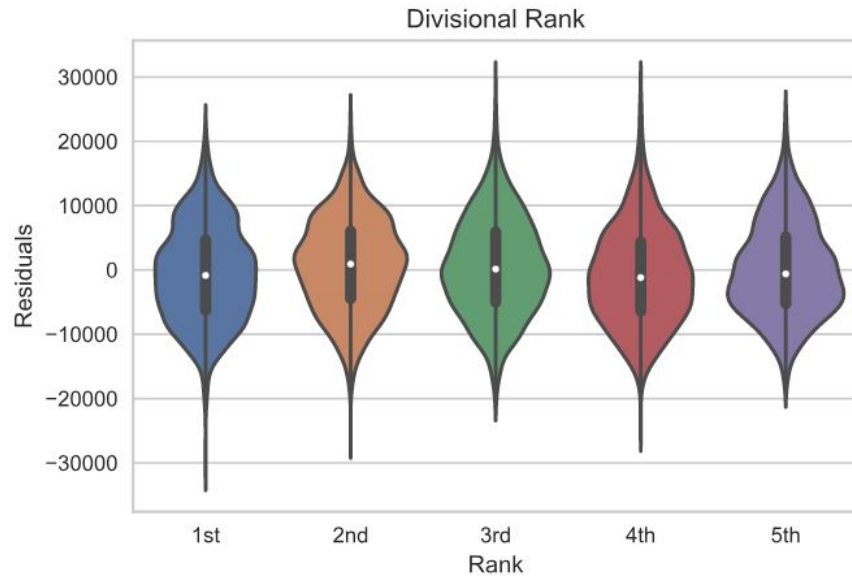
Lasso (Lamba \approx 2.05, Standardized)

- Rank
- Games back
- Night or day
- Wins last 10
- Run differential
- Win differential
- Mean batter age
- # of current all-stars
- # of lifetime all-stars
- Team salary
- Day of week

RMSE \approx 7434.25

$R^2 \approx 0.44$

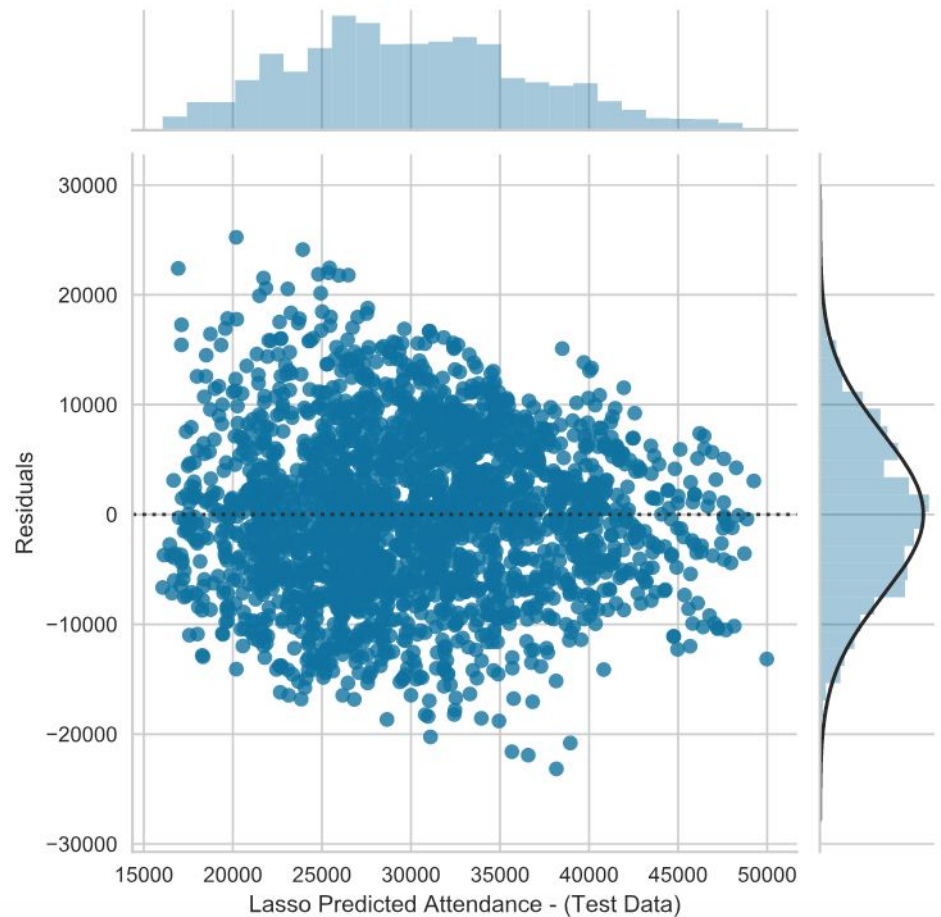
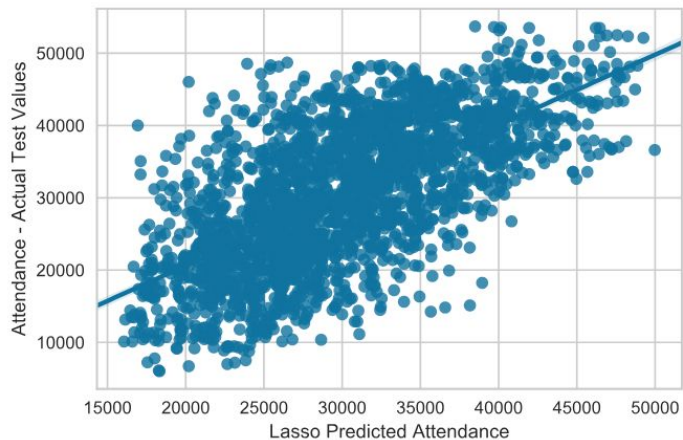
Checking Residuals



*Training Data

Outcome

- Lasso Model
 - RMSE ≈ 7124.2
 - $R^2 \approx 0.434$



Future Work

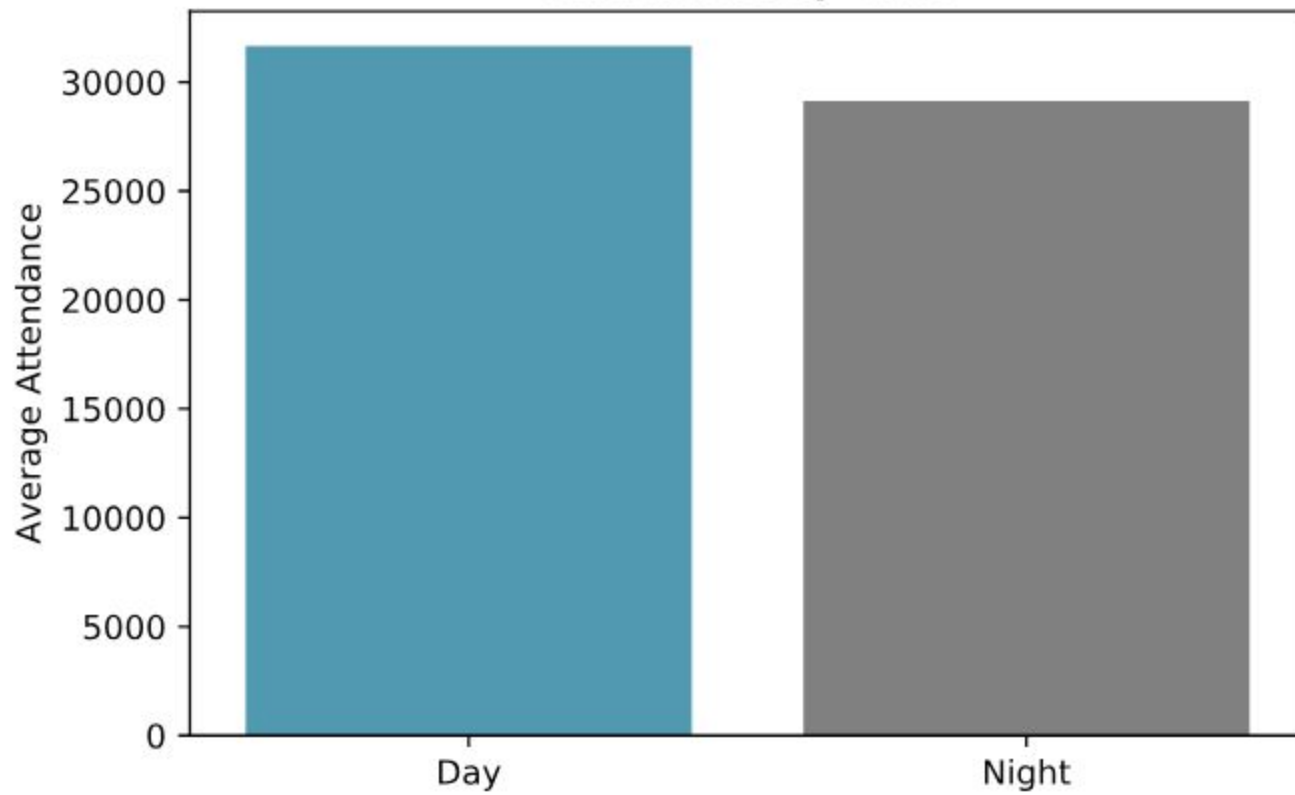


- Transforming Y variable
- Adding more features
 - Weather
 - Digging in further on opponent features
 - Other sports leagues

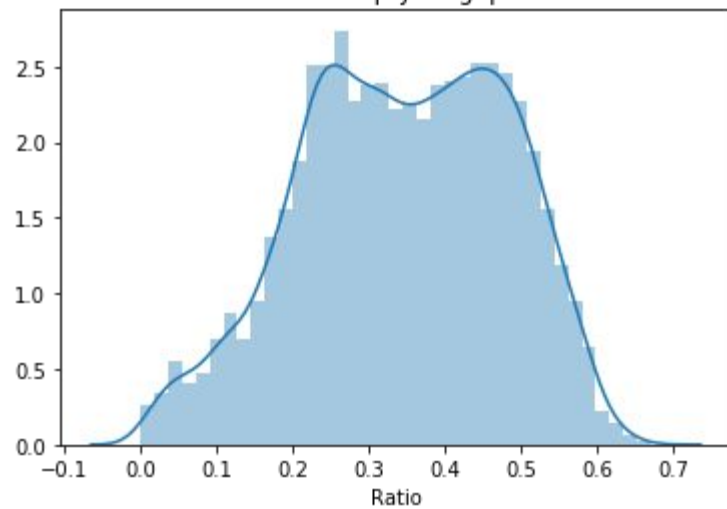
Questions?

Appendix

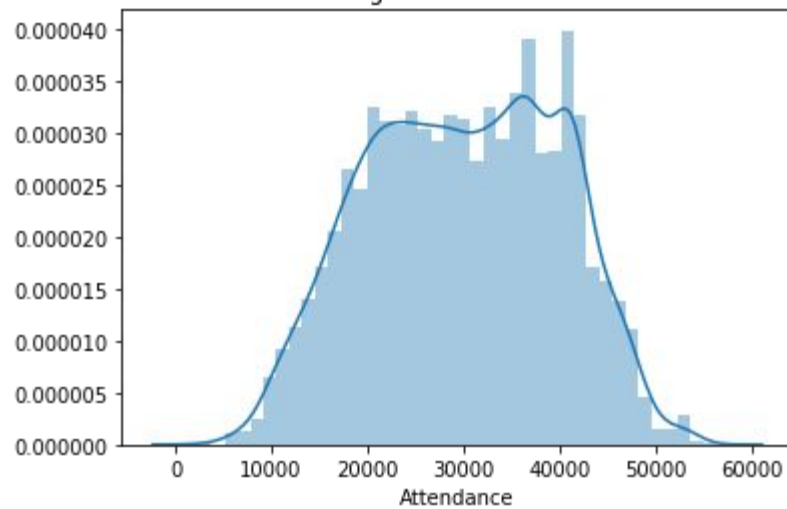
Attendance by Time

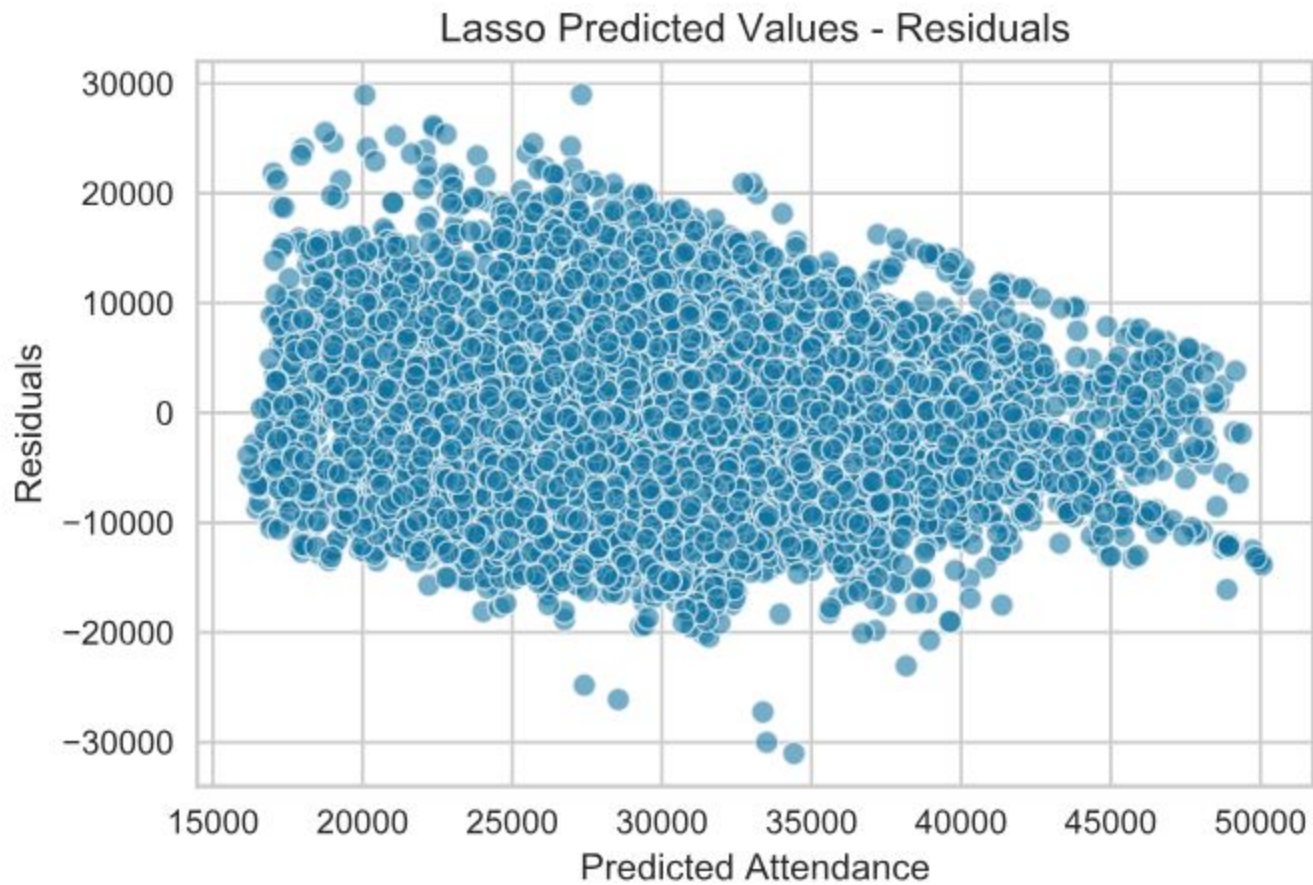


Ratio of Stadium Empty - log1p Transformed

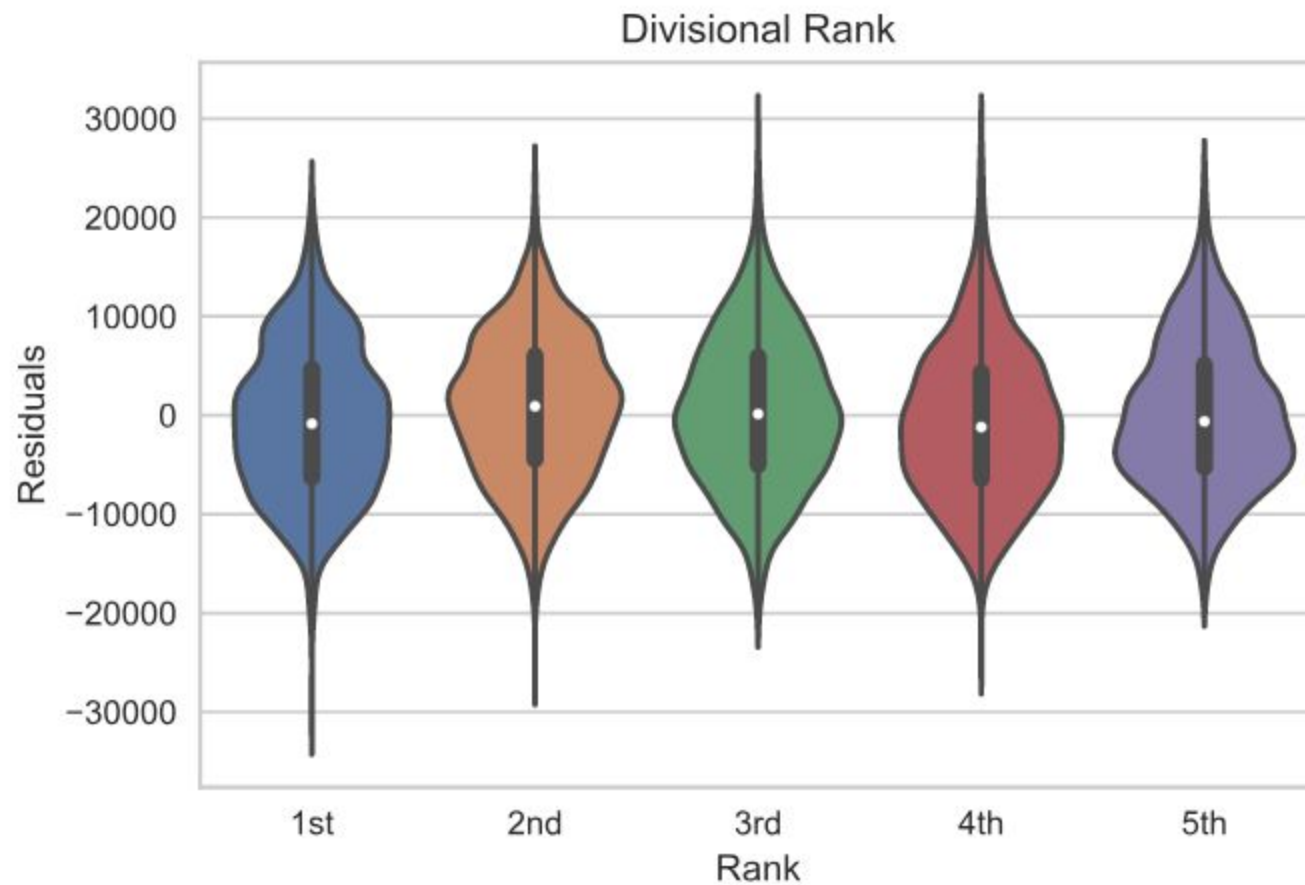


Regular Attendance

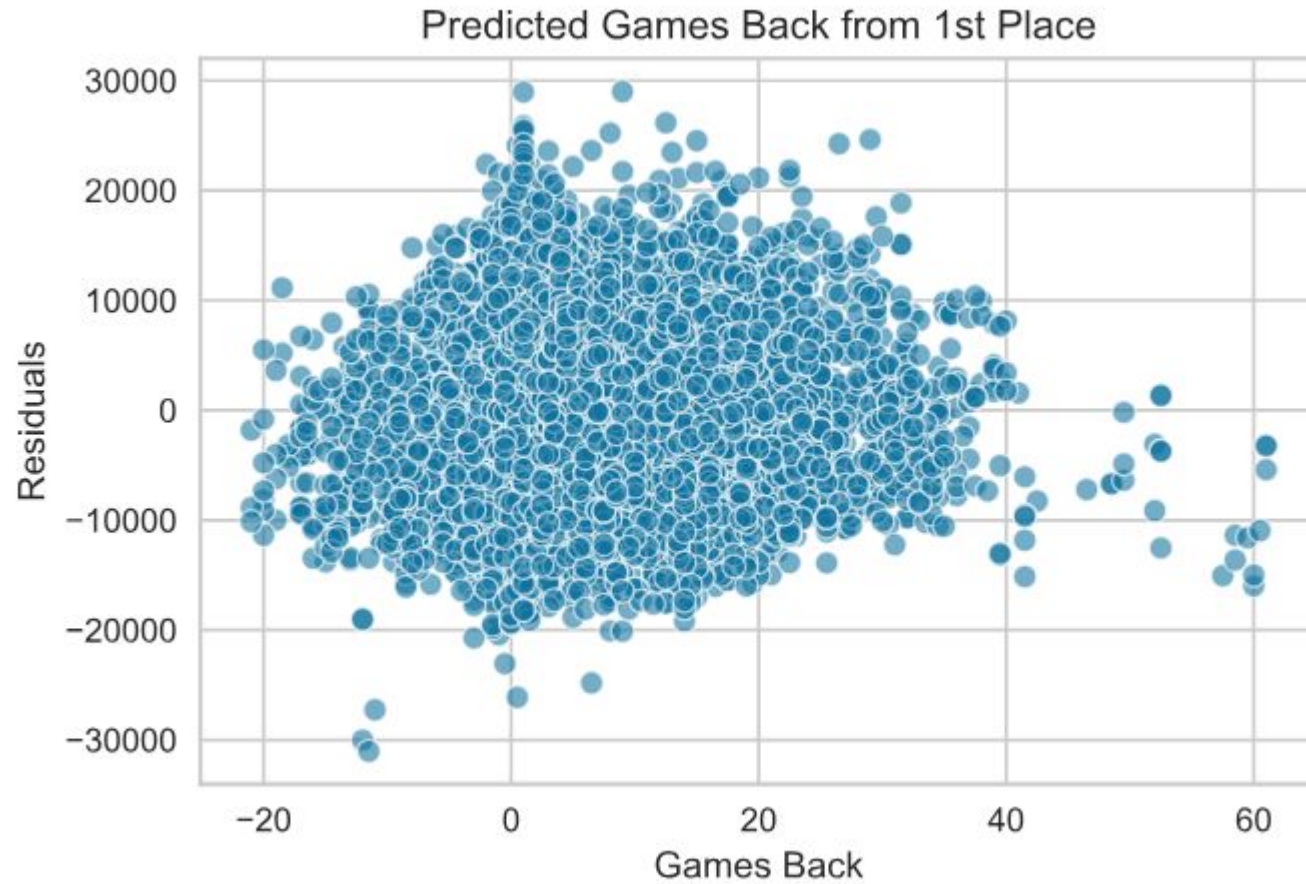




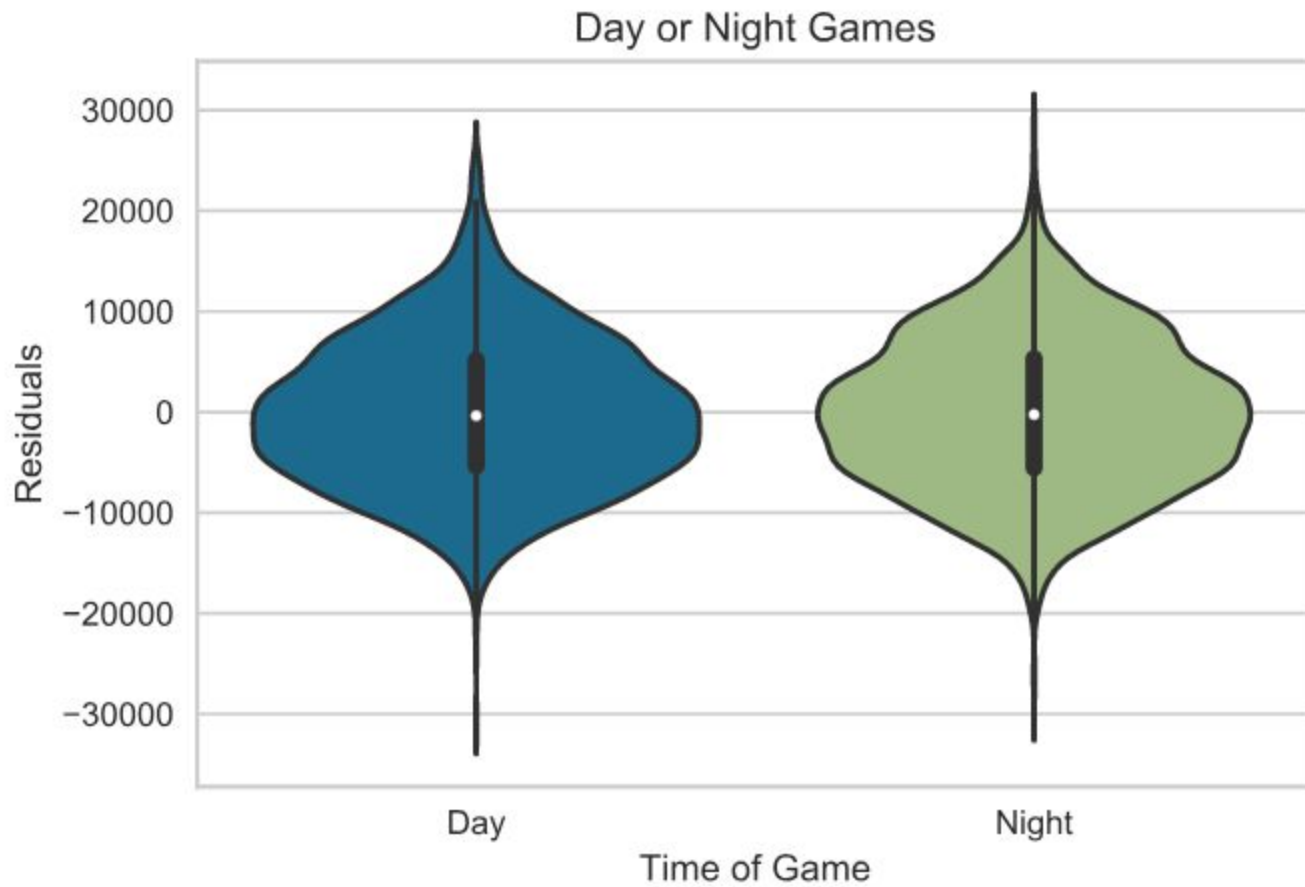
*Training Data



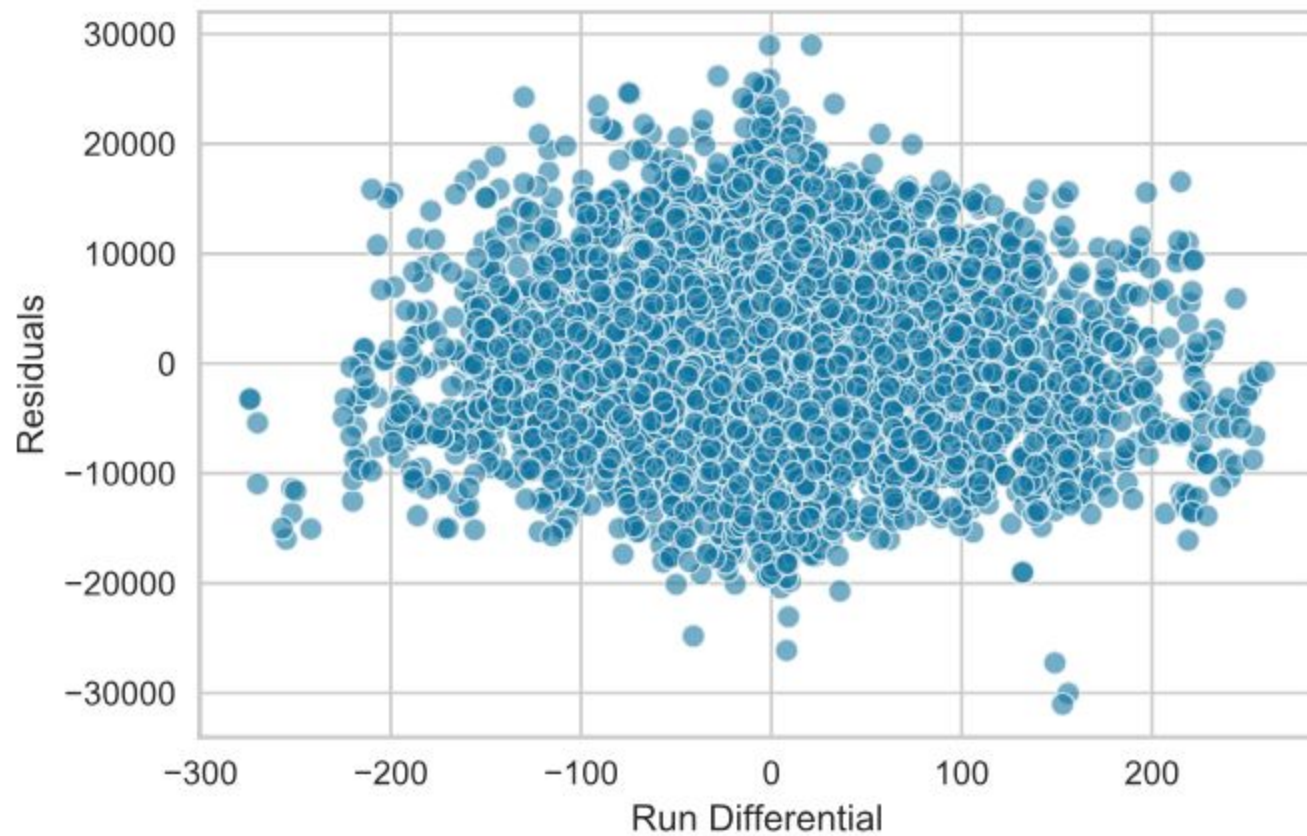
*Training Data



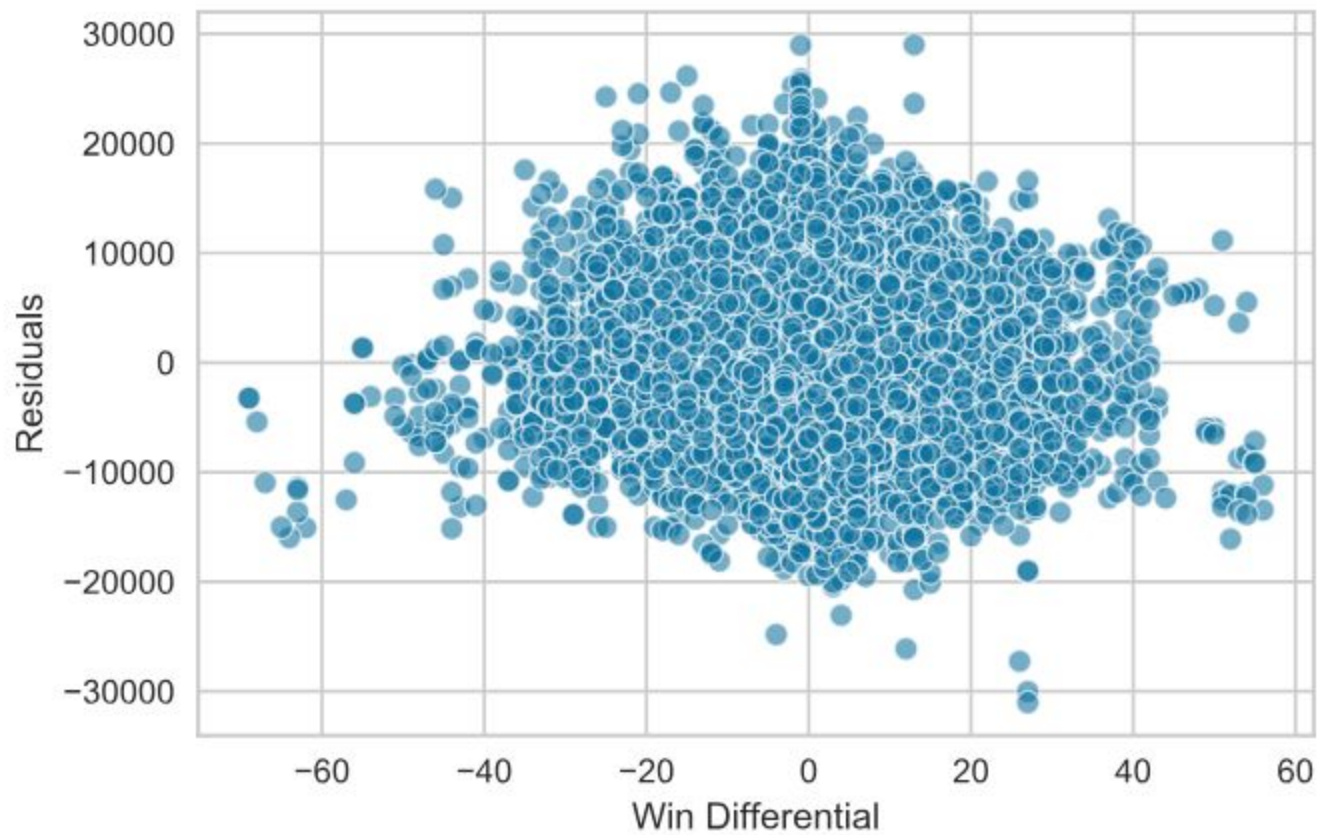
*Training Data

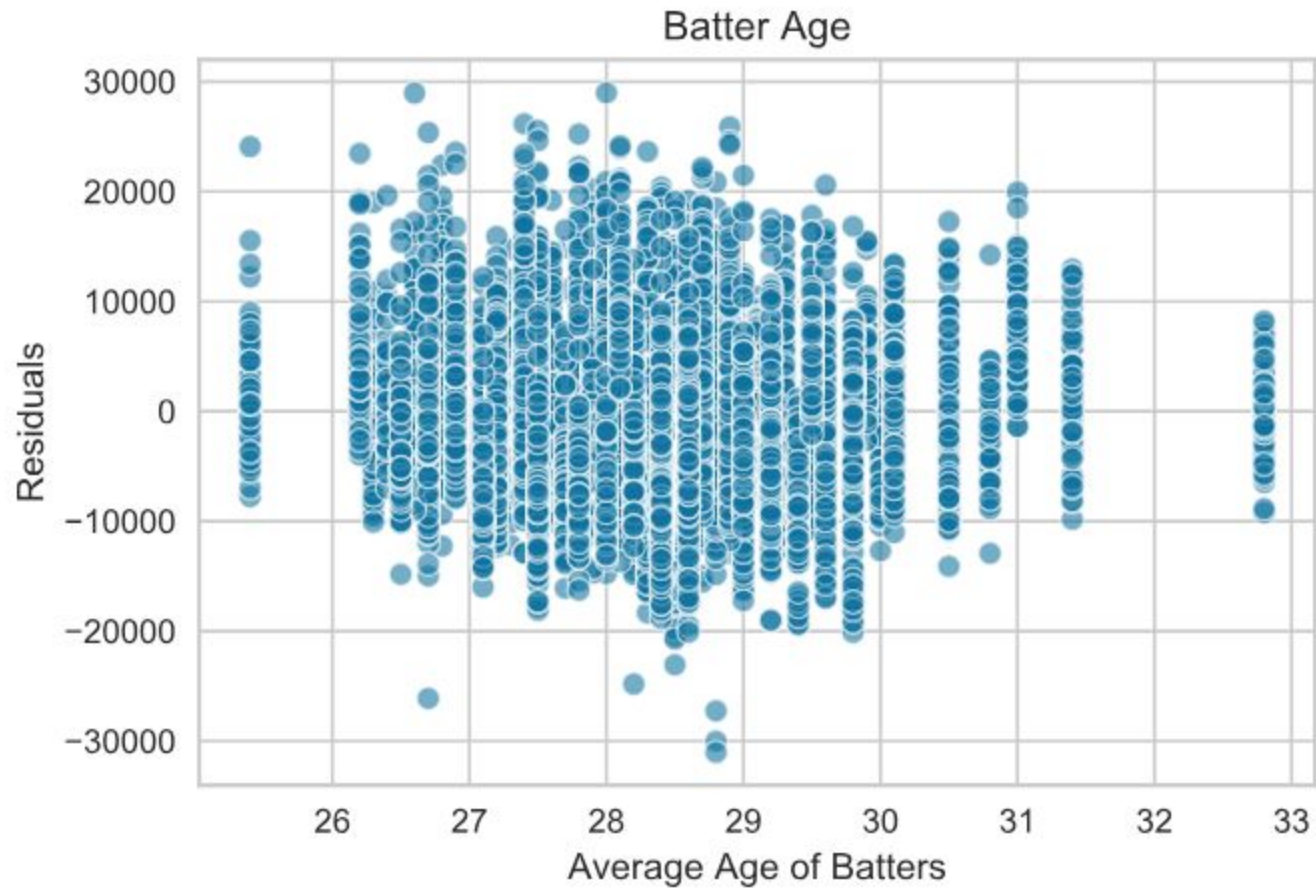


*Training Data

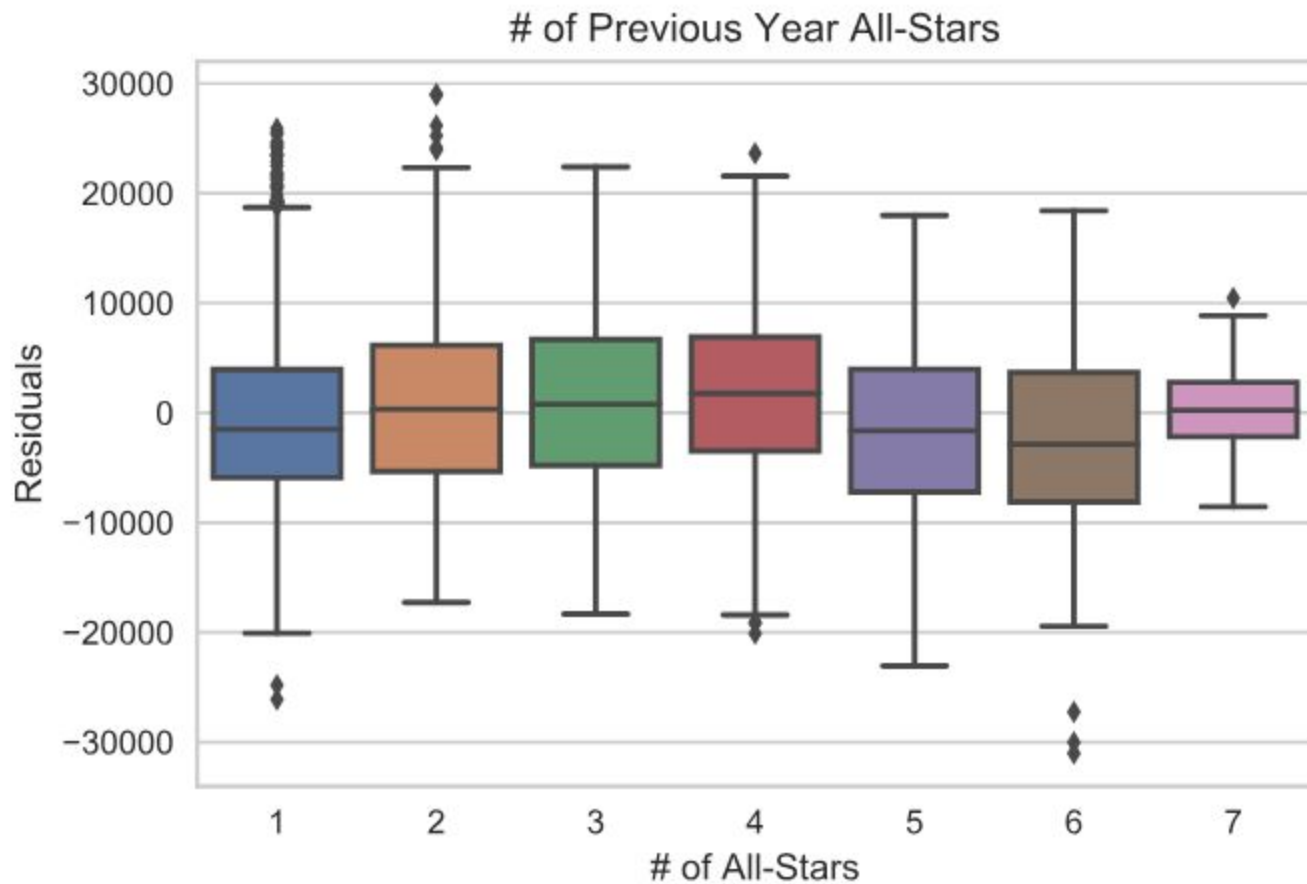


*Training Data

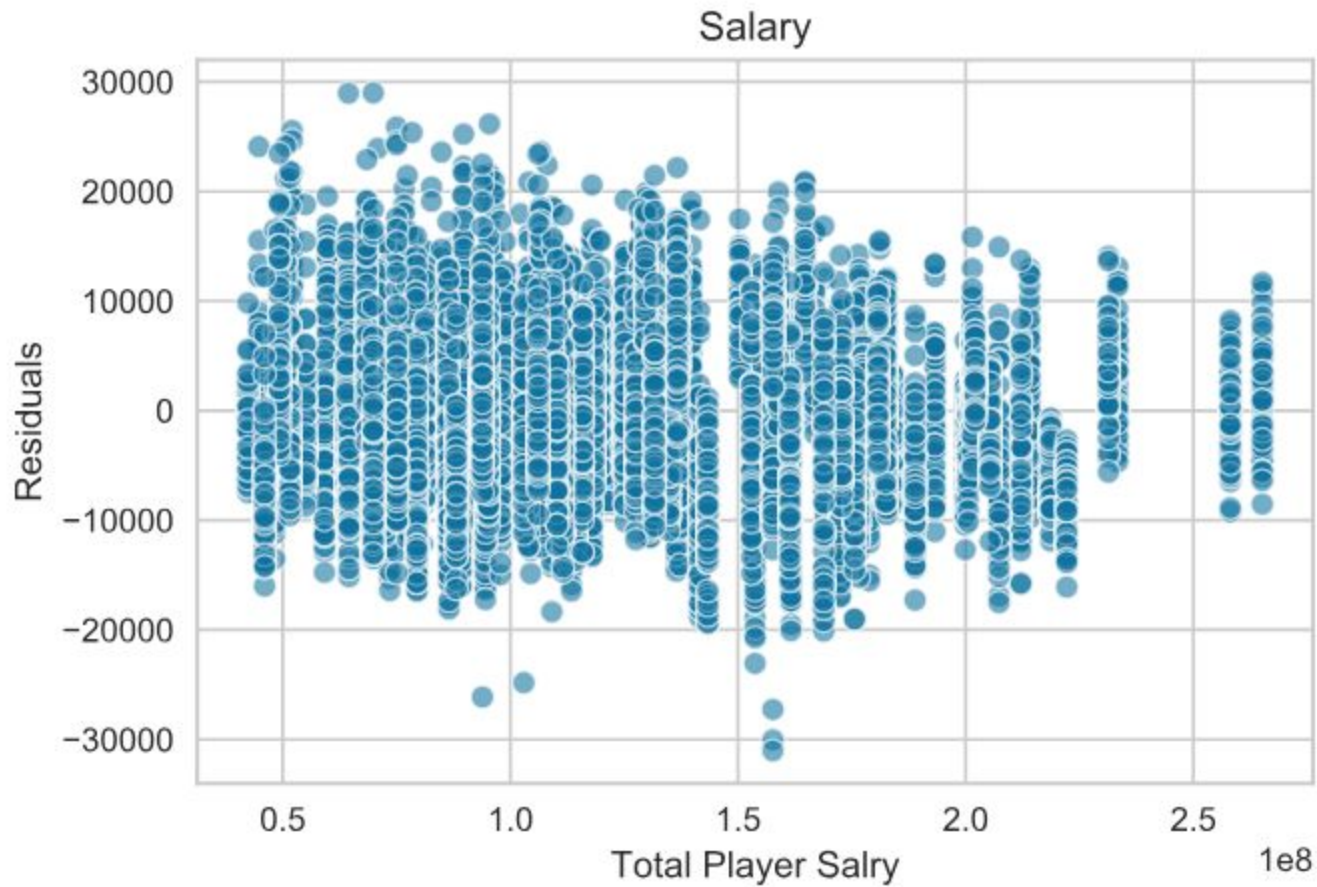




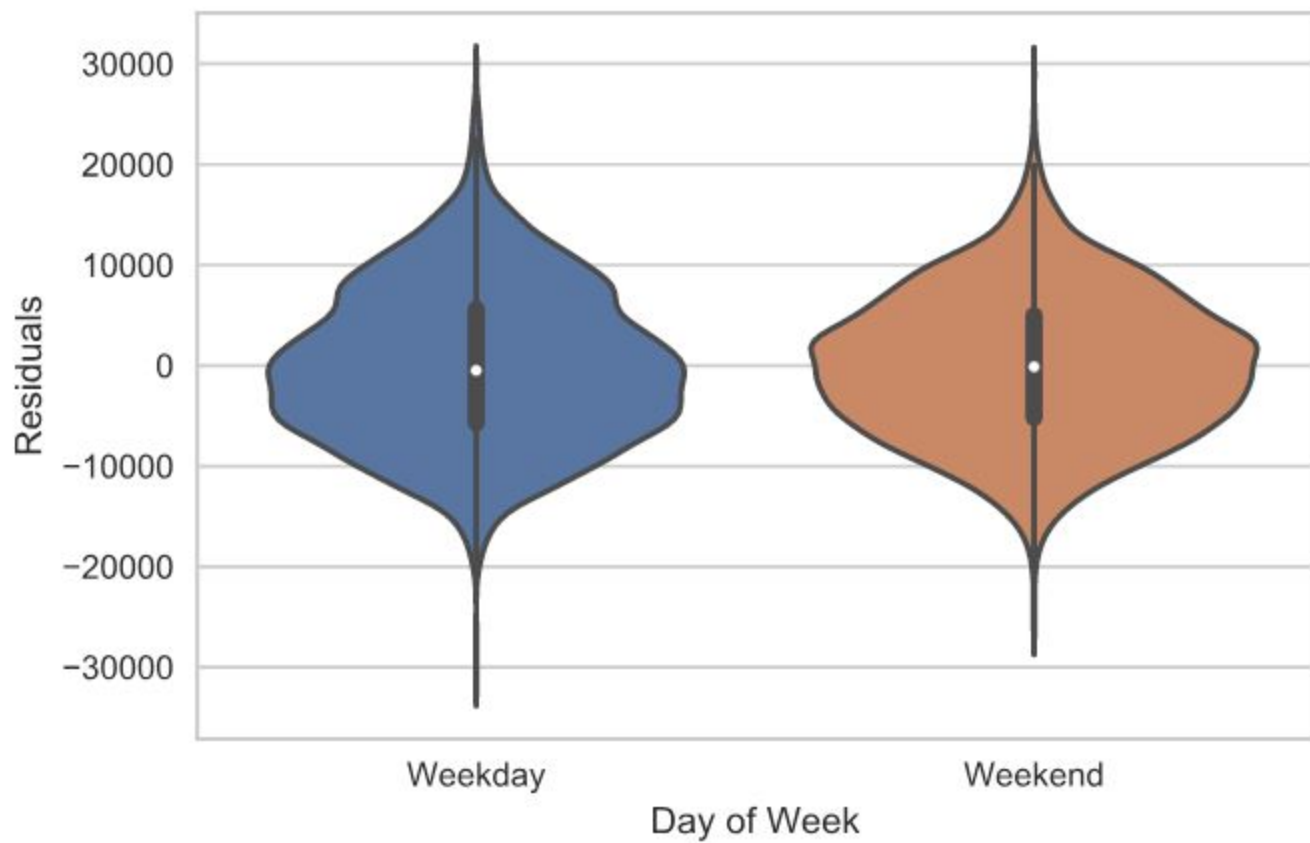
*Training Data



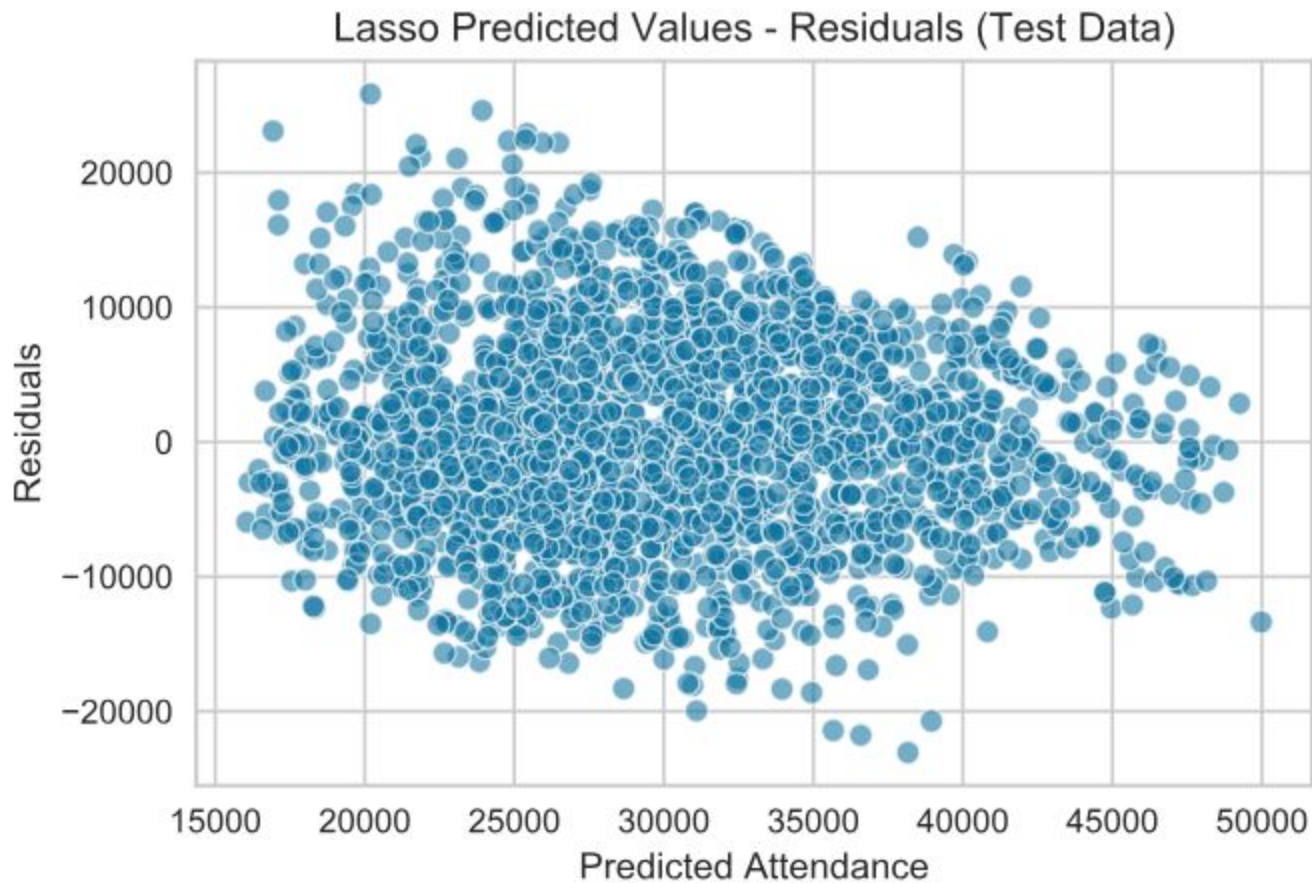
*Training Data



*Training Data



*Training Data



*Test Data

Opposing Team Features - Heat Map

