# How Efficient Baserunning Impacts Scoring

By: Henry Bednar

## Executive Summary

This comprehensive analysis employed advanced machine learning techniques to develop predictive baserunning models for two critical offensive performance metrics: Run Scoring Percentage (RSP) and Runs Per Game (RPG). Through rigorous feature engineering, complete coefficient analysis, and model validation processes, it successfully identified the key performance drivers with statistical precision.

### Key Findings:

- Achieved robust predictive models with R² values of 0.586 for RSP and 0.534 for RPG
- Identified 8 statistically significant predictors for RSP and 9 statistically significant predictors for RPG (p < 0.05)
- Quantified impact measurements including -5.093 coefficient for speed (RSP) and -2.738 for speed (RPG)
- Discovered counter-intuitive insights such as negative coefficients for stolen base opportunities
- Established temporal stability with year coefficients proving statistically irrelevant (p > 0.05 for both models)

## Research Methodology & Enhanced Analytical Framework

### Advanced Analytical Approach

This study employed a comprehensive machine learning pipeline enhanced with complete coefficient extraction and statistical significance testing. The enhanced framework incorporated:

- **Complete Coefficient Analysis**: Linear regression coefficient extraction with standard errors, p-values, and 95% confidence intervals for all features
- **Statistical Significance Framework**: Systematic p-value analysis identifying reliable predictors from total features

- **Multi-Algorithm Ensemble**: Stacked ensemble achieving 1.59% improvement for RSP (though individual models outperformed for RPG)
- **Rigorous Cross-Validation**: 5-fold cross-validation ensuring consistent performance across data configurations
- **Advanced Feature Selection**: Boruta algorithmic selection (RSP: 35 features, RPG: 32 features)

## Data Scope & Model Performance

The analysis encompasses Baseball Savant and Baseball Reference data from 270 team-seasons across nine complete seasons (2016-2024):

- **RSP Model Performance**: $R^2$ = 0.586, Adjusted $R^2$ = 0.508, F-statistic = 7.52 (35 features, 181 DF)
- **RPG Model Performance**: $R^2$ = 0.534, Adjusted $R^2$ = 0.452, F-statistic = 6.54 (32 features, 183 DF)
- **Temporal Validation**: Both models show stable relationships across all nine seasons

# Coefficient Analysis - Run Scoring Percentage (RSP)

## Model Equation

- RSP = -40.335 + 0.133×(On_2nd_Single_Scored) + 0.009×(On_1st_When_Double)
-    - 5.093×(Average_Time_Home_to_First) - 0.005×(Stolen_Base_Opportunities)
-    + 0.131×(On_1st_Double_Scored) + ... (29 more terms)

### Tier 1: Statistically Significant Predictors (p < 0.05)

**On_2nd_Single_Scored**

- Coefficient: +0.133 | p-value: 0.014* | 95% CI: [0.027, 0.239]
- Each 1% increase in scoring from 2nd base on singles adds 0.133 percentage points to overall team scoring efficiency—the ultimate clutch metric

**Stolen_Base_Opportunities**

- Coefficient: -0.005 | p-value: < 0.001*** | 95% CI: [-0.007, -0.003]
- Counter-intuitive finding: Excessive stealing opportunities may indicate offensive struggles—teams that create many steal chances but don't convert them into runs through hitting

**Bases_Taken**

- Coefficient: +0.016 | p-value: 0.008** | 95% CI: [0.004, 0.028]
- Aggressive baserunning philosophy—each additional base taken per opportunity adds 0.016 percentage points to scoring efficiency

## On_2nd_Single_Reached_3rd

- Coefficient: +0.100 | p-value: 0.035* | 95% CI: [0.007, 0.192]
- Critical RISP advancement—moving from 2nd to 3rd on singles positions runners in prime scoring position

## Number_of_Players

- Coefficient: -0.088 | p-value: 0.012* | 95% CI: [-0.157, -0.019]
- Roster instability indicator—teams using more players typically struggle with offensive consistency and chemistry

## CS_Only_Runner_On

- Coefficient: -0.146 | p-value: 0.014* | 95% CI: [-0.261, -0.030]
- Pure baserunning mistakes when the runner is the only one on base—highly costly situational errors

## Outs_On_Base_2nd

- Coefficient: -0.075 | p-value: 0.022* | 95% CI: [-0.139, -0.011]
- Rally-killing mistakes at second base—each additional out costs 0.075 percentage points in scoring efficiency

## SB_2B_Advances_vs_Average

- Coefficient: +0.072 | p-value: 0.036* | 95% CI: [0.005, 0.139]
- Above-average advancement quality on stolen bases—not just stealing, but stealing effectively

## Tier 2: High-Impact Non-Significant Contributors

### On_1st_When_Double

- Coefficient: +0.009 | p-value: 0.918 | Importance: 0.547
- Baseball IQ metric with highest importance—ensuring that runners are on 1st when doubles occur maximizes scoring potential

### Average_Time_Home_to_First

- Coefficient: -5.093 | p-value: 0.132 | Importance: 0.485

- Speed foundation with massive coefficient—each 0.1 second improvement could add ~0.5 percentage points to scoring

**On_1st_Double_Scored**

- Coefficient: +0.131 | p-value: 0.138 | Importance: 0.453
- Aggressive advancement from 1st to home on doubles—substantial potential impact when executed

# Coefficient Analysis - Runs Per Game (RPG)

## Model Equation

- RPG = 95.077 - 0.0004×(Stolen_Base_Opportunities) - 2.738×(Average_Time_Home_to_First)
- + 0.010×(On_1st_When_Double) + 0.003×(Bases_Taken)
- + 0.280×(Lead_Distance_Pitcher_Windup_to_Release) + ... (27 more terms)

### Tier 1: Statistically Significant Predictors (p < 0.05)

**Average_Time_Home_to_First**

- Coefficient: -2.738 | p-value: 0.001** | 95% CI: [-4.398, -1.077]
- Speed foundation for run production—each 0.1 second improvement adds ~0.27 runs per game over a full season

**On_2nd_Single_Scored**

- Coefficient: +0.030 | p-value: 0.009** | 95% CI: [0.008, 0.052]
- RISP conversion drives run volume—each percentage point improvement adds 0.030 runs per game

**Lead_Distance_Pitcher_Windup_to_Release**

- Coefficient: +0.280 | p-value: 0.015* | 95% CI: [0.055, 0.504]
- Technical baserunning skill—optimal lead distances create multiple scoring opportunities per game

**Outs_On_Base_2nd**

- Coefficient: -0.021 | p-value: 0.024* | 95% CI: [-0.039, -0.003]
- Rally-killing mistakes at second base with measurable per-game cost

### On_2nd_When_Single

- Coefficient: -0.025 | p-value: 0.023* | 95% CI: [-0.046, -0.004]
- Counter-intuitive negative—may indicate forcing advancement rather than smart situational baserunning

### CS_Only_Runner_On

- Coefficient: -0.040 | p-value: 0.021* | 95% CI: [-0.074, -0.006]
- Baserunning mistakes with quantified run cost when no other runners are present

### Number_of_Players

- Coefficient: -0.036 | p-value: < 0.001*** | 95% CI: [-0.053, -0.019]
- Roster instability strongly indicates offensive struggles—highly significant across both models

### Average_Sprint_Speed

- Coefficient: -0.297 | p-value: 0.029* | 95% CI: [-0.563, -0.030]
- Quality over pure speed—suggests that context-dependent application matters more than raw athleticism

### On_2nd_Single_Reached_3rd

- Coefficient: +0.025 | p-value: 0.027* | 95% CI: [0.003, 0.048]
- Consistent contributor to aggressive advancement from scoring position

## Tier 2: High-Impact Non-Significant Contributors

### Stolen_Base_Opportunities

- Coefficient: -0.0004 | p-value: 0.096 | Importance: 0.667
- Too many stealing opportunities may indicate an inability to drive runners in

### On_1st_When_Double

- Coefficient: +0.010 | p-value: 0.599 | Importance: 0.343
- Baseball IQ metric with practical but statistically variable impact

### Bases_Taken

- Coefficient: +0.003 | p-value: 0.102 | Importance: 0.287
- Aggressive baserunning philosophy with borderline statistical significance

# Cross-Model Strategic Insights

### Statistically Validated Universal Success Factors

**Scoring from 2nd on Singles**

- RSP Impact: +0.133 coefficient (p = 0.014*) - clutch efficiency validated
- RPG Impact: +0.030 coefficient (p = 0.009**) - volume production confirmed
- Strategic Value: Ultimate clutch metric validated across both efficiency and volume systems

**Baserunning Discipline (Multiple "Getting Out" Metrics)**

- Both Models: Consistent negative coefficients with statistical significance
- RSP Impact: Outs at 2nd (-0.075, p = 0.022*), CS only runner (-0.146, p = 0.014*)
- RPG Impact: Outs at 2nd (-0.021, p = 0.024*), CS only runner (-0.040, p = 0.021*)

**Roster Management Discipline**

- RSP Impact: Number of players used (-0.088, p = 0.012*)
- RPG Impact: Number of players used (-0.036, p < 0.001***)
- Strategic Value: Roster stability correlates strongly with offensive success

### Counter-Intuitive Statistical Discoveries

**Stolen Base Opportunities Paradox**

- Both Models: Negative coefficients (-0.005 RSP, -0.0004 RPG)
- Statistical Significance: RSP highly significant (p < 0.001***), RPG marginal (p = 0.096)
- Strategic Implication: Excessive stealing opportunities likely indicate inability to drive runners in throughl hitting

**Speed Application Complexity**

- Average Time Home to First: Large negative coefficients (speed improvement beneficial)
- Average Sprint Speed: Negative RPG coefficient (-0.297, p = 0.029*) suggests context matters more than raw speed
- Strategic Value: Intelligent speed application trumps pure athleticism

## Temporal Stability Analysis

### Year Effect Statistical Analysis

**RSP Model Year Impact**

- Coefficient: +0.049 | p-value: 0.644 | Improvement: 0.25% RMSE
- Interpretation: Statistically irrelevant temporal trend

**RPG Model Year Impact**

- Coefficient: -0.035 | p-value: 0.180 | Improvement: -0.79% RMSE
- Interpretation: Year actually reduces model performance—ignore temporal adjustments

**Strategic Implication**: Fundamental baserunning and situational skills transcend rule changes and analytical evolution.

# Statistical Significance Hierarchy

## Highly Significant Predictors (p < 0.01)

**RSP Model:**

- Stolen Base Opportunities: -0.005 coefficient (p < 0.001***)
- Bases Taken: +0.016 coefficient (p = 0.008**)

**RPG Model:**

- Average Time Home to First: -2.738 coefficient (p = 0.001**)
- On_2nd_Single_Scored: +0.030 coefficient (p = 0.009**)
- Number of Players Used: -0.036 coefficient (p < 0.001***)

# Summary of Statistical Findings

## RSP Model:

- Total features analyzed: 35
- Statistically significant features (p < 0.05): 8 out of 35
- Highest coefficient impact: Average_Time_Home_to_First (-5.093, non-significant)
- Most reliable predictor: Stolen_Base_Opportunities (p < 0.001)

## RPG Model:

- Total features analyzed: 32
- Statistically significant features (p < 0.05): 9 out of 32
- Highest coefficient impact: Average_Time_Home_to_First (-2.738, p = 0.001)
- Most reliable predictor: Number_of_Players (p < 0.001)

# Key Strategic Takeaways

1. **Clutch Advancement**: Moving runners from 2nd to 3rd and scoring from 2nd on singles are universally valuable across both efficiency and volume metrics

2. **Discipline Over Aggression**: Avoiding outs on the basepaths, especially at 2nd base and when caught stealing as the only runner, provides measurable value

3. **Speed Context Matters**: Raw speed is less important than intelligent application—optimal leads and situational awareness trump pure athleticism

4. **Roster Stability**: Teams using fewer players throughout the season demonstrate better offensive consistency and baserunning execution

5. **Quality Over Quantity**: Excessive stolen base opportunities may indicate offensive deficiencies rather than aggressive baserunning prowess