

Evaluating Hitting Skills of NPB Players with Logistic Regression Analysis

Mamoru Sakai¹, Hiroki Tanioka², Kenji Matsuura², Masahiko Sano², Kenji Ohira²,
Tetsushi Ueta² and Hiroaki Sakaguchi³

¹Graduate School of Advanced Technology and Science
Tokushima University, Japan

²Center for Administration of Information Technology
Tokushima University, Japan

³Shikoku Island League Plus, Japan

Abstract—*Luck or unluck is supposed to exist in baseball. In this paper, we consider unpredictable “lucky” or “unlucky” cases which are occurred on the ground in order to find out authentic ability of batters. If all cases are observable, hit probabilities can be calculated. However, it is impossible to gather all the information on the ground. Therefore, “luck” is defined as influence of information that cannot be observed. To provide more appropriate evaluation for NPB¹ players, this paper proposes a predicted BABIP using Logistic Regression and describes the predicted BABIP reducing the influence of luck on the assessment of hitting ability. Analysis results show that there is a possibility of appropriate evaluation on hitting abilities of NPB players using the predicted BABIP.*

Keywords: baseball, hitting, performance, evaluation, logistic regression

1. Introduction

For the professional baseball league in Japan, statistical data has been indispensable to evaluate players and make a tactical plan. A baseball scorebook existed more than 100 years ago. The scorebook has a history of being utilized for making strategy and tactical plan of baseball game. Currently, a large amount of data which means a big data is easily collected as information technology advances. There is a possibility of eliminating subjectivity with the big data. That is why every professional baseball team applied a variety of evaluation systems using the big data. On the other hand, every baseball player needs an appropriate evaluation because the evaluation result strongly influences to the team plan and the baseball player's career as an athlete.

1.1 Background of Research

When we watch a baseball game, we often witness scenes that a line drive flies out in front of a fielder. On the other hand, a pop fly falls between a position of infielders and

¹NPB (Nippon Professional Baseball) is a professional baseball league in Japan, which is often called as “Japanese Baseball”.

outfielders. In such a scene, we viewers express “unlucky” or “lucky” for the batter. Therefore, luck or unluck is supposed to be existing in baseball. In other words, it is hard to say that hitting records depend on only the ability of the batters. Conversely, assuming that all factors including position of a ball hit, quality of a ball, position of fielding, ability of fielding, and shape of a stadium, in this case, the hitting result must be predicted completely. However, there are unpredictable cases using the limited information from the baseball game. In this paper, we consider that unpredictable cases are viewed as “lucky” or “unlucky” by human eyes. Hence, luck is defined as the influence by information that an observer cannot observe.

1.2 Purpose of Research

Player's performance is appeared based on the player's abilities affected by luck. A case study [1] said that the randomness affects player's hitting streaky. When the player is compelled to play under unlucky situation, unfortunately, no matter how good play, it will be just treated as a failure. If lucky players are overestimated, they can perform under their estimated abilities in the next season. Contrariwise, if unlucky players are underestimated, they can perform over their estimated abilities in the next season. Thus, players and their team owners need an appropriate indicator of performance which produces their stable results. In other words, the indicator describes player's authentic ability. Therefore, a purpose of this research is to clarify lucky players and unlucky players by using an indicator which is based on differences between results of two seasons in 2015 and 2016. The data including those results is provided by a workshop [2] under the sponsorship of The Institute of Statistical Mathematics.

2. Indicators of baseball players

2.1 Sabermetrics

There are many Key Performance Indicators (KPIs) of baseball players, such as Batting Average (AVG), the number of Home Runs (HRs), Runs Batted In (RBI), etc. These days,

various methods and indicators are proposed to evaluate performance of players more objectively. One of typical indicators is Sabermetrics (Society for American Baseball Research Metrics) [3], [4]. Sabermetrics was proposed in the 1970s by Bill James who is the baseball writer. Sabermetrics is an objective method for analyzing baseball data and evaluating players. Major League Baseball (MLB) official records are based on Sabermetrics.

2.2 Operations Research

Operations Research [5] is a strategy for evaluating baseball players [6]. It is one of methodology to help leaders make better decisions using mathematical and statistical models.

2.3 BABIP

There is a problem in the Sabermetrics and related researches, those which are not considering where a play is happened in the ground. BABIP (Batting Average On Ball In Play) was proposed by Voros McCracken [7][8]. BABIP is the percentage of hits on ground except Home Runs in batted ball. The BABIP equation is:

$$BABIP = \frac{H - HR}{AB - SO - HR + SF}, \quad (1)$$

where H is Hits, HR is Home Runs, AB is At Bats, SO is Strikeouts, and SF is Sacrifice Flies.

Sasaki reported the result [9] that there was low correlation between consecutive two seasons in BABIP of NPB players. According to Chris Dutton's research using Regression Analysis [10], the batting result does not depend only on the ability of the batter, but the opponent's defense and the ground environment. Thus, an indicator should be considered in consideration of the situations on the ground to evaluate true abilities of the players.

3. Proposed Method

In this work, batting data is used as learning data, and a regression model is created using logistic regression analysis with a target variable. The target variable is a binary variable whether the batting was a hit or not. Next, a predicted BABIP as a theoretical value is calculated by the created regression model with observable information at the bat. Lastly, every player is evaluated by comparing the predicted BABIP and an actual BABIP.

3.1 Linear Regression

Chris Dutton designed a regression model to determine the relationship between each factor and a hitter's BABIP [10]. Quantitative variables can be used for explanatory variables of linear regression. From a batting data, the hitting information is set as explanatory variables, and the batting result

whether hit or not (1 or 0) is set as a target variable.

$$p = b_0 + \sum_{j=1}^k b_j x_j, \quad (2)$$

where p is a predicted value based on x , x_1, \dots, x_k are explanatory variables, b_0 is a constant, and b_1, \dots, b_k are regression coefficients.

3.2 Logistic Regression

Logistic Regression Analysis [11] is employed for regression modeling in order to obtain hit probability of each batting. Both quantitative variables and qualitative variables can be used for explanatory variables of logistic regression. However, the regression model doesn't include the situation on the ground. In our research, the hitting information including the situation on the ground from batting data is set as explanatory variables, and the batting result whether hit or not (1 or 0) is set as a target variable.

As explanatory variables are set on Eq. (3) as regression model, the hit probability is calculated as a value $[0, 1]$. When a batted ball except Home Runs flows into the ground, the hit probability is within the range of 0 to 1. Here, Strikeout is as 0 and Home Run is as 1. The equation of logistic regression model is given as follows.

$$\log \frac{p}{1-p} = b_0 + \sum_{j=1}^k b_j x_j, \quad (3)$$

where p is a predicted value based on x , x_1, \dots, x_k are explanatory variables, b_0 is a constant, and b_1, \dots, b_k are regression coefficients. Transformed equation for p is expressed as follows.

$$p = \frac{1}{1 + \exp \left(-b_0 - \sum_{j=1}^k b_j x_j \right)}. \quad (4)$$

3.3 Explanatory Variable

Explanatory variables are chosen among only the information on the ground after hitting. As the ground is regarded as a lottery box, the ground information influences the distribution of lottery tickets. The following information is adopted as explanatory variable.

- 1) Coordinates of grounder the ball
Coordinates (x, y) of grounder, fly and line drive are converted to polar coordinates (r, θ) .
- 2) Runner situation of each base
If a runner is on a base, the runner situation is set as 1, otherwise set as 0.
- 3) Defense strength of the opponent team
DER (Defense Efficiency Rating) [12] is an indicator of the team's defense strength.

Table 1: DER (Defense Efficiency Rating) of 12 NPB teams in 2015 and 2016.

year	Giants	Tigers	Dragons	Carp	Baystars	Swallows	Hawks	Lions	Fighters	Buffaloes	Marines	Eagles
2015	0.712	0.678	0.696	0.691	0.674	0.696	0.703	0.701	0.693	0.691	0.677	0.676
2016	0.680	0.687	0.696	0.700	0.688	0.686	0.707	0.674	0.703	0.678	0.685	0.668

Table 2: An example of explanatory variables for regression model.

grounder r	grounder θ	fly r	fly θ	line drive r	line drive θ	1st base	2nd base	3rd base	DER	Hits
0.00	0.00	178.00	0.45	0.00	0.00	0	0	0	0.712	0
0.00	0.00	213.80	1.27	0.00	0.00	0	0	0	0.712	0
161.12	1.29	0.00	0.00	0.00	0.00	0	0	0	0.674	1
110.26	0.90	0.00	0.00	0.00	0.00	1	0	0	0.674	0
0.00	0.00	222.32	1.24	0.00	0.00	0	1	0	0.674	1
0.00	0.00	30.61	0.90	0.00	0.00	0	1	0	0.674	0
0.00	0.00	181.37	0.85	0.00	0.00	0	1	0	0.674	0
88.53	0.84	0.00	0.00	0.00	0.00	0	0	0	0.712	0
115.80	1.28	0.00	0.00	0.00	0.00	0	0	0	0.712	0
106.02	1.13	0.00	0.00	0.00	0.00	0	0	0	0.712	0
0.00	0.00	159.13	0.29	0.00	0.00	1	1	0	0.712	0
81.99	1.05	0.00	0.00	0.00	0.00	1	1	0	0.712	0

* DER is Defense Efficiency Rating of opponent team.

Table 3: Result of Linear Regression Analysis.

	Estimate	std.Error	z value	p value
constant	0.1048	0.0878	1.194	0.2327
grounder r	0.0065	4.622×10^{-5}	141.203	2.00×10^{-16}
grounder θ	0.0199	0.0046	4.295	1.75×10^{-5}
fly r	0.0041	0.0003	109.785	2.00×10^{-16}
fly θ	0.0301	0.0041	7.646	2.10×10^{-16}
line drive r	0.0078	5.838×10^{-5}	133.134	2.00×10^{-16}
line drive θ	-0.0269	0.0093	-2.915	0.0036
First base	0.0188	0.0033	5.693	1.26×10^{-8}
Second base	0.0076	0.0038	1.988	0.0468
Third base	0.0443	0.0052	8.570	2.00×10^{-16}
DER	-0.8382	0.1271	-6.592	4.37×10^{-11}

* DER is Defense Efficiency Rating of opponent team.

Table 4: Result of Logistic Regression Analysis.

	Estimate	std.Error	z value	p value
constant	-3.2198	0.6661	-4.833	1.34×10^{-6}
grounder r	0.0497	0.0004	102.809	2.00×10^{-16}
grounder θ	0.4937	0.0404	12.206	2.00×10^{-16}
fly r	0.0319	0.0003	90.296	2.00×10^{-16}
fly θ	0.5523	0.0346	15.946	2.00×10^{-16}
line drive r	0.0617	0.0008	76.025	2.00×10^{-16}
line drive θ	-0.4749	0.1048	-4.532	5.85×10^{-6}
First base	0.0751	0.0249	3.009	0.0026
Second base	0.0591	0.0286	2.065	0.0389
Third base	0.3216	0.0391	8.221	2.00×10^{-16}
DER	-6.0978	0.9644	-6.323	2.57×10^{-10}

* DER is Defense Efficiency Rating of opponent team.

It is difficult to calculate defense strength to each player. Instead, DER is employed for an explanatory variable.

$$DER = \frac{PA - H - BB - HBP - SO - E}{PA - HR - BB - HBP - SO}, \quad (5)$$

where PA is Plate Appearances, H is Hits, BB is Bases on Balls, HBP is Hit by Pitch, SO is Strikeouts, and E is Errors in the numerator. Then, PA is Plate Appearances, HR is Home Runs, BB is Bases on Balls, HBP is Hit by Pitch, and SO is Strikeouts in the denominator.

Table 1 shows DER of 12 teams in 2015 and 2016, respectively. Table 2 shows an example of explanatory variables used for regression analysis.

3.4 Evaluation method of players in luck

The actual BABIP (ABA) is the player's practical BABIP throughout the season. The predicted BABIP (PBA) is an expectation value based on hit probabilities of a player. A luck score ($Luck$) is the actual BABIP and a difference between the predicted BABIP.

$$Luck = ABA - PBA. \quad (6)$$

It is good luck if the actual BABIP is higher than the predicted BABIP, and it is bad luck if the actual BABIP is lower. The players over 200 plate appearances for each

season are targeted for luck analysis. There are 100 players in 2015 and 109 players in 2016. To give players an appropriate evaluation, the luck score should be deducted in BABIP. If the predicted BABIP excluding luck is stable between consecutive two seasons, the predicted BABIP is an appropriate indicator.

4. Experiment

4.1 Regression Model Analysis

Table 3 is output results of the linear regression analysis, which shows estimated coefficients and related statistics. Table 4 is output results of the logistic regression analysis, which shows estimated coefficients and related statistics.

The equation of linear regression model in which the coefficient of the explanatory variable is substituted into the Eq. (2) is given as follows.

$$p = 0.1048 + 0.0065x_1 + \dots - 0.8382x_{10} \quad (7)$$

The equation of logistic regression model in which the coefficient of the explanatory variable is substituted into the Eq. (4) is given as follows.

$$p = \frac{1}{1 + \exp(3.2198 - 0.0497x_1 - \dots + 6.0978x_{10})} \quad (8)$$

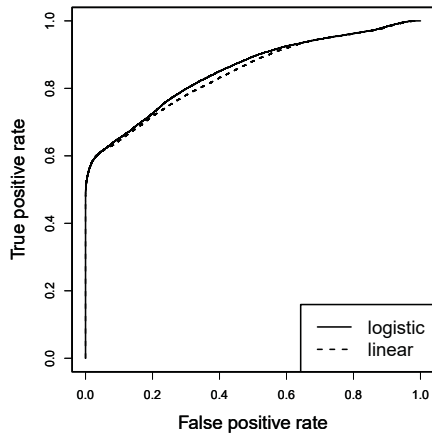


Fig. 1: ROC (Receiver Operating Characteristic) curve.

Table 5: Actual and Prediction results.

		Prediction results	
		Out	Hit
Actual results	Out	59,452	1,173
	Hit	11,220	15,483

Table 6: Statistics and Accuracies (cutoff value is 0.59)

statistics	values	statistics	values
True Positive	11,500	Sensitivity	0.57
False Positive	682	Precision	0.94
False Negative	8,556	Specificity	0.98
True Negative	42,761	Accuracy	0.85

4.2 Prediction of batting averages

Whether a hit probability goes to hit or not, that depend on a cutoff value. The cutoff value could be [0, 1]. Both true positive rate and false positive rate are moved depending on the cutoff value. The accuracy of a test is its ability to differentiate actual hitting results and predicted hitting results correctly.

Figure 1 shows two curves which are drawn for the linear regression and logistic regression. The curve is called as the ROC (Receiver Operating Characteristic) curve. The logistic regression curve is upper than the linear regression curve. This means that the logistic regression is better as an indicator than the linear regression. The cutoff value is chosen on the ROC curve following *Accuracy*.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \quad (9)$$

where *Accuracy* is the accuracy, *TP* is true positive, *TN* is true negative, *FP* is false positive, and *FN* is false negative. Table 6 shows the cutoff value, the result of discrimination and each accuracy index from Table 5. Where a cutoff value is 0.59, an accuracy is 0.85 which is the maximum value. Table 5 shows a relationship between actual and prediction results as the cutoff is 0.59.

Table 7: Statistics of each distribution

	Counts	Average	Variance	SD
Actual BABIP (2015)	100	0.312	0.001372	0.037
Actual BABIP (2016)	109	0.315	0.001343	0.0366
Predicted BABIP (2015)	100	0.316	0.000845	0.0291
Predicted BABIP (2016)	109	0.311	0.00079	0.0281
Luck (2015)	100	-0.00444	0.001008	0.0317
Luck (2016)	109	0.00389	0.000882	0.0297

Table 8: Lucky players in 2015

	name	actual	predict	luck	L	R	speed
1	Atsushi Fujii	0.385	0.310	0.075	1	1	4.66
2	Takehiro Ishikawa	0.345	0.273	0.072	1	0	4.86
3	Soichiro Tateoka	0.372	0.306	0.066	1	0	5.89
4	Haruki Nishikawa	0.348	0.284	0.064	1	0	7.76
5	Shingo Kawabata	0.377	0.319	0.059	1	0	3.28
6	Takuya Nakajima	0.328	0.269	0.059	1	0	5.58
7	Ryota Imanari	0.397	0.347	0.050	1	0	2.93
8	Ryo Hijirisawa	0.346	0.297	0.049	1	0	5.73
9	Tsuyoshi Ueda	0.307	0.266	0.041	1	0	5.04
10	Kyohei Kamezawa	0.316	0.275	0.041	1	0	4.70

Table 9: Unlucky players in 2015

	name	actual	predict	luck	L	R	speed
1	Miguel Mejia	0.289	0.385	-0.096	0	1	0.99
2	Brandon J. Laird	0.244	0.310	-0.096	1	0	2.09
3	Shuichi Murata	0.264	0.316	-0.052	1	0	1.57
4	Tsubasa Aizawa	0.286	0.334	-0.048	1	0	2.79
5	Tatsuihiro Tamura	0.216	0.264	-0.048	1	0	3.94
6	Luis Cruz	0.271	0.318	-0.047	1	0	2.34
7	Keiji Obiki	0.262	0.308	-0.046	1	0	3.85
8	Mitsutaka Goto	0.252	0.298	-0.046	0	1	3.78
9	Seiichi Uchikawa	0.303	0.348	-0.045	1	0	1.96
10	Hiroyuki Nakajima	0.287	0.331	-0.044	1	0	1.88

4.3 Actual BABIP and Predicted BABIP

Table 7 shows the statistics of each distribution. Figure 3 and Figure 4 show the actual batting distribution and the predicted batting distribution. T-test was conducted to determine whether there was a significant difference between the actual BABIP and the predicted BABIP. The null hypothesis states that "There is no difference between the actual BABIP and the predicted BABIP". As a result, *p* is 0.349 (≥ 0.05 , $df = 198$) in 2015, also *p* is 0.382 (≥ 0.05 , $df = 216$) in 2016. Thus, the null hypothesis could not be rejected, and there is not significant difference.

4.4 Lucky players and Unlucky players

Table 8 shows the top ten players whose actual BABIP was higher than the predicted BABIP in 2015. And Table 9 shows top ten players whose actual BABIP was lower than the predicted BABIP in 2015. According to the Eq. (6), lucky players increased their actual BABIP due to the luck score, and unlucky players decreased their actual BABIP due to the luck score. In other words, it is a player who are overestimated and underestimated in 2015. Additionally, Table 8 includes L, R, and speed. L and R mean left-handed and right-handed respectively, and speed is the speed score in Eq. (10).

Table 13: Luck distributions in two seasons 2015–2016.

	mean	variance	players	df	p value
Luck Dist	-9.57×10^{-5}	0.000964	209	208	-
Luck Dist L/R	0.000282	0.000877	209	208	0.247
Luck Dist Spd	0.000617	0.000734	209	208	0.0249

players looks smaller in the L/R values and the speed scores. In quantitative observation, the luck score distributions are compared in Figure 5 and Figure 6. Each distribution follows a normal distribution. Table 13 shows statistical information including variance values of the luck score distributions in 2015–2016.

There is not a statistically significant reduction 0.000087 between Luck Dist and Luck Dist L/R in variance. There is a statistically significant reduction 0.000023 between Luck Dist and Luck Dist Spd in variance. Here, p is 0.0249 (< 0.05) in F-test. From these results, the predicted BABIP got closer to the actual BABIP, while the luck score distribution was shrunk. This means that the more observable information is adopted to logistic regression as explanatory variable, the more the luck score distribution is shrunk. If all the information on the ground is observable, luck is not existing in baseball games.

6. Conclusion

In this research, on an assumption that there is luck that the players cannot control in baseball. We proposed an indicator “predicted BABIP” using logistic regression to evaluate players properly. The predicted BABIP as a theoretical value was calculated by the created regression model with observable information at the bat. Every player was evaluated by the predicted BABIP. Assuming that a difference between the actual BABIP and the predicted BABIP is a luck score, lucky players and unlucky players are separated by the luck scores.

Investigating the tendency of lucky players and unlucky players, lucky players hit many grounders and could be fast runners, unlucky players hit many fly balls and might be slow runners. It seems that there are still factors to be considered such as player’s ability in the part which was influenced by luck. In addition, it seems that the influence which the running ability was not considered. When the L/R values and the speed score were adopted to logistic regression as explanatory variables, the bias derived from player’s ability was slightly decreased in the luck score.

It was found from the result that the predicted BABIP properly evaluated player’s hitting ability rather than the actual BABIP, because the predicted BABIP has smaller influence of luck than the actual BABIP. However, it is not enough to gathering observable information on the ground,

because, a speed score is a pseudo speed score which does not represent a player’s running ability to the first base every at bat.

Therefore, we need to measure every arrival time for reaching the first base after hitting. Then, the future works arising from this research are to clarify what is the player’s running ability, to eliminate all the luck on the ground, and to provide more appropriate indicator to evaluate every baseball players.

Acknowledgment

This research was supported by Data Stadium Inc. and The Institute of Statistical Mathematics. We thank our colleagues from A6 Computer System Laboratory of Tokushima University who provided insight and expertise that greatly assisted the research, although they may not agree with all the conclusions of this paper. We thank Michitomo Morii for assistance with particular technique, and Kohei Kawanaka for comments that greatly improved the manuscript.

References

- [1] S. C. Albright, “A Statistical Analysis of Hitting Streaks in Baseball,” *Journal of the American Statistical Association*, vol. 88, no. 424, pp. 1175–1183, 1993.
- [2] M. Sakai and H. Tanioka, “Yakyu ni Okeru Un toha Nanika - Logistic Kaiki Bunseki wo Mochiita Anda Kakuritsu no Yosoku (What is Lucy in Baseball? – Prediction of Hit Probability Using Logistic Regression),” The Institute of Statistical Mathematics, Tech. Rep., March, published in Japanese, 2018.
- [3] M. M. Lewis., “Moneyball: The Art of Winning an Unfair Game,” NY, USA, 2003.
- [4] R. J. Puerzer, “From Scientific Baseball to Sabermetrics: Professional Baseball as a Reflection of Engineering and Management in Society,” *NINE: A Journal of Baseball History and Culture*, vol. 11, no. 1, pp. 34–48, 2002, accessed: 2018-02-14.
- [5] C. W. Churchman, R. L. Ackoff, Ackoff, and E. Arnoff, “Introduction to operations research,” 1957.
- [6] N. Streib, S. J. Young, and J. Sokol, “A Major League Baseball Team Uses Operations Research to Improve Draft Preparation,” vol. 42, pp. 119–130, March 2012.
- [7] V. McCracken, “Pitching and Defense: How Much Control Do Hurlers Have?” <https://www.baseballprospectus.com/news/article/878/pitching-and-defense-how-much-control-do-hurlers-have/>, January 2001, accessed: 2018-02-14.
- [8] D. Studeman, “Data Erratum Et Cetera,” <https://www.fangraphs.com/tht/data-erratum-et-cetera/>, January 2004, accessed: 2018-02-14.
- [9] H. Sasaki, “BABIP ga Imisuru Tokoro to, sono Kaishaku no Muzukashisa (The meaning of BABIP and the Difficulty of the Interpretation),” <http://www.baseball-lab.jp/column/entry/175/>, May, published in Japanese, 2015, Accessed: 2018-02-14.
- [10] C. Dutton, “Batters and BABIP,” <https://www.fangraphs.com/tht/batters-and-babip/>, December 2008, accessed: 2018-02-14.
- [11] S. R. Bailey, “Forecasting Batting Averages in MLB,” Master’s thesis, Simon Fraser University, November 2017.
- [12] “Defensive Efficiency Ratio (DER),” <http://m.mlb.com/glossary/advanced-stats/defensive-efficiency-ratio>, accessed: 2018-02-14.
- [13] B. James, *The Bill James Baseball Abstract 1987 (1st ed.)*. Ballantine Books, 1987, accessed: 2018-02-14.

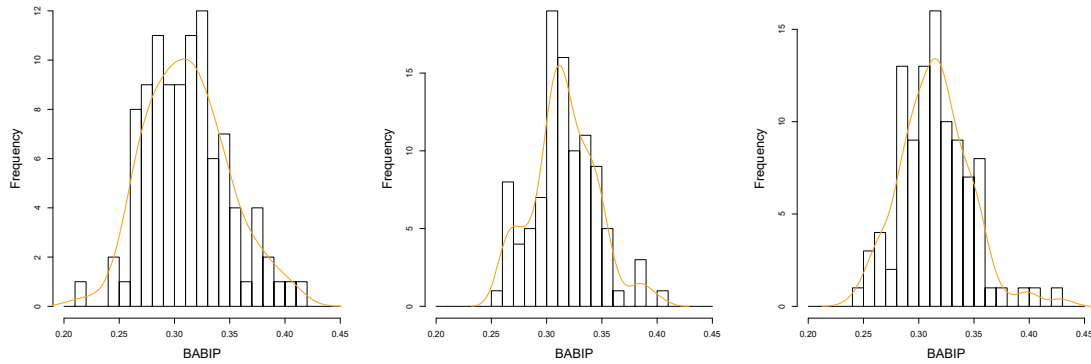


Fig. 3: Distributions of BABIP, predicted BABIP, and predicted BABIP with speed score in 2015.

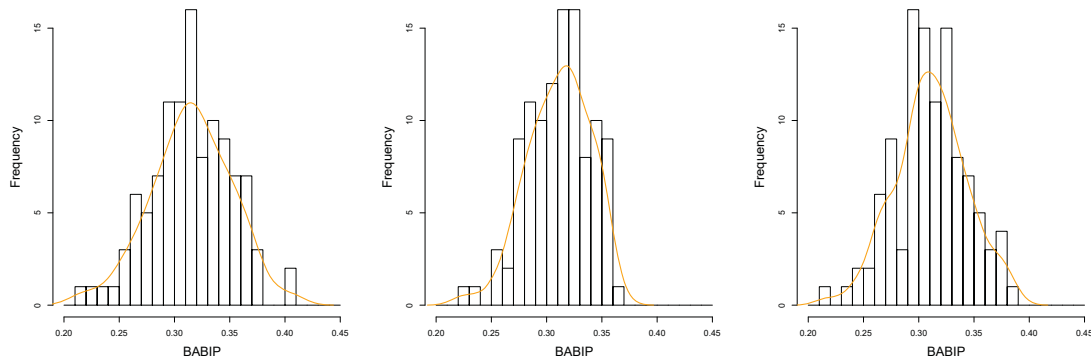


Fig. 4: Distributions of BABIP, predicted BABIP, and predicted BABIP with speed score in 2016.

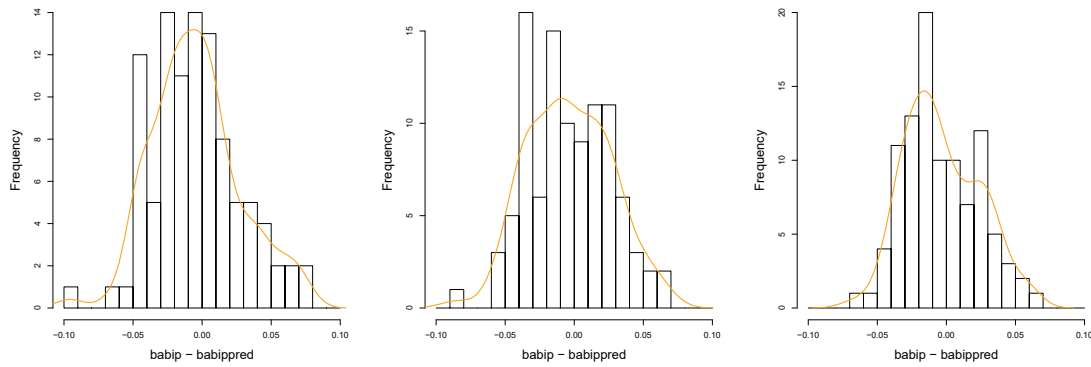


Fig. 5: Distributions of luck, luck with L/R, and luck with speed score in 2015.

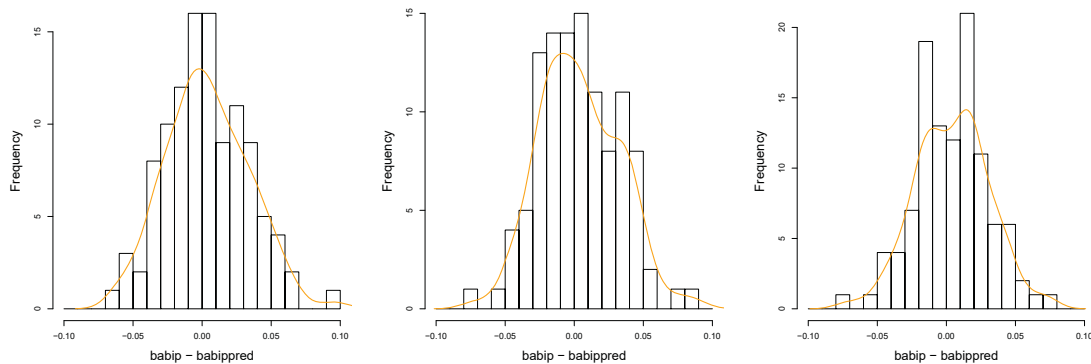


Fig. 6: Distributions of luck, luck with L/R, and luck with speed score in 2016.