

Evaluating plate discipline in Major League Baseball with Bayesian Additive Regression Trees

Ryan Yee* Sameer K. Deshpande†

May 11, 2023

We introduce a three-step framework to determine, on a per-pitch basis, whether batters in Major League Baseball should swing at a pitch. Unlike traditional plate discipline metrics, which implicitly assume that all batters should always swing (resp. take) pitches inside (resp. outside) the strike zone, our approach explicitly accounts not only for the players and umpires involved but also in-game contextual information like the number of outs, the count, baserunners, and score. Specifically, we first fit flexible Bayesian nonparametric models to estimate (i) the probability that the pitch is called a strike if the batter takes the pitch; (ii) the probability that the batter makes contact if he swings; and (iii) the number of runs the batting team is expected to score following each pitch outcome (e.g. swing and miss, take a called strike, etc.). We then combine these intermediate estimates to determine whether swinging increases the batting team’s run expectancy. Our approach enables natural uncertainty propagation so that we can not only determine the optimal swing/take decision but also quantify our confidence in that decision. We illustrate our framework using a case study of pitches faced by Mike Trout in 2019.

1 Introduction

1.1 Motivating example

At the top of the seventh inning of the September 5, 2019 game between the Oakland Athletics and the Los Angeles Angels, Angels batter Kevan Smith faced Athletics pitcher A.J. Puk with the Angels leading 5 runs to 1. On an 0-2 pitch with no outs and no runners on-base, Puk threw a fastball that just missed the lower right-hand corner of the strike zone (see Figure 1). Smith decided to swing at the pitch. Did Smith make the correct decision?

*Department of Statistics, University of Wisconsin–Madison. ryee2@wisc.edu

†Department of Statistics, University of Wisconsin–Madison. sameer.deshpande@wisc.edu

By swinging, Smith risked missing the pitch, picking up a strike, and losing an out, thereby putting his team at a disadvantage. And even if he had made contact, he risked flying or grounding out, which would similarly disadvantage his team. At the same time, however, by swinging, Smith had a chance of getting on base or even scoring a home run. On the other hand, had Smith not swung and had instead taken the pitch, because the pitch just missed the inside edge of the strike zone, the umpire may have correctly called the pitch a ball (which would advantage Smith’s team) or mistakenly called a strike (which would disadvantage Smith’s team).



Figure 1: Kevan Smith swings at a fastball thrown just outside the lower right-hand corner of the strike zone.

As it turns out, Smith hit a home run on this pitch. So at least retrospectively, it would seem like the decision to swing was good. But was Smith simply lucky? Or should we expect him to hit home runs consistently off of similar pitches in similar game situations? And how sure should we be about these expectations?

We provide quantitative answers to these questions using a Bayesian modeling framework to assess batter’s decision making — or plate discipline — prior to observing the outcome of each pitch. In the case of Smith in Figure 1, we find that there was a 12.3% chance (90% credible interval [5.8%, 16.7%]) that the umpire would call the pitch a strike. Although taking the pitch would likely have benefitted Smith’s team, Smith had about an equal chance of making contact if he swung (83.1%, 90% credible interval [80.6%, 85.3%]). In fact, we find

that by swinging on similar pitches, Smith increases the number of runs the Angels are expected to score in the remainder of the half-inning only 61.4% of the time and only by about 0.01 runs on average (90% credible interval $[-0.04, 0.08]$). Ultimately, our analysis concludes that while Smith’s decision to swing on the pitch was not especially risky, it was unlikely to substantially benefit his team, in terms of the number of expected runs the team would score in the rest of the inning.

Traditional plate discipline metrics compare the proportion of pitches batters swing at that are outside and inside the strike zone (Slowinski, 2010a). In this way, these metrics reward batters for avoiding taken strikes based on the implicit assumption that pitches thrown outside (resp. inside) the strike zone are always called balls (resp. strikes). Of course, umpires do not adhere strictly to the official strike zone when making ball/strike decisions. Though they often rely on adaptive heuristics and prior experience to make calls (Green and Daniels, 2022), umpires can also be influenced by the framing ability of individual catchers (Marchi, 2011; Lindbergh, 2013; Deshpande and Wyner, 2017); player age or ability (Kim and King, 2014; Mills, 2014); their previous calls (Chen et al., 2016); and the fact that their calls are reviewed by the league (Mills, 2017). As a result, traditional plate discipline metrics may systematically penalize batters who swing at “frameable” pitches that just barely miss the strike zone but are likely to be called strikes. These metrics additionally fail to account for the fact that some pitches are easier to hit than others. At least intuitively, we might regard a batter who only takes hard-to-hit pitches as more disciplined than a batter who consistently takes easy-to-hit pitches. Further, traditional plate discipline metrics entirely ignore the many contextual and situational factors that can influence a batter’s decision to swing. For instance, batters may be more aggressive on pitches thrown with two strikes to avoid striking out while looking. Finally, existing plate discipline metrics do not quantify the uncertainty in their findings.

1.2 Our contributions

We introduce a three-step approach to determine optimal swing decisions, assess batter decision-making on a per-pitch basis, and quantify our uncertainty about both. In the first step, we fit flexible Bayesian nonparametric models that enable us to estimate, for any single pitch, (i) the probability that a batter makes contact; (ii) the probability that the umpire calls a strike; and (iii) the average number of runs that the batting team is expected to score in the remaining half-inning as functions of the pitch location, players and umpires involved,

and game-state information like the count, inning, and baserunners. Then, in the second step, we combine these estimates using the law of total expectation to compute the expected number of runs the batting team is expected to score if the batter swings or takes the pitch. Finally, we use these quantities to determine the optimal swing/take decision. By fitting Bayesian models in the first step, we are able to propagate uncertainty about the actual outcomes of a particular pitch through to our assessment of the batter’s decision making in a natural and computationally efficient manner.

Our framework involves fitting three component models, each of which may be of independent interest. First, elaborating on work begun in [Deshpande and Wyner \(2017\)](#), we develop a Bayesian model of called strike probability that “borrows statistical strength” across umpires and players through flexible partial pooling of data. We demonstrate that this model, which is based on Bayesian additive regression trees (BART; [Chipman et al., 2010](#)) and accounts not only for pitch location but also player and umpire identities and game-state information, outperforms parametric competitors based on generalized additive models that account only for location. We develop a similar model for the probability that a batter makes contact on a pitch. Finally, we developed a BART-based run expectancy model, which we call **BARTxR**, to predict the number of runs a team is expected to score following each pitch outcome as a function of in-game contextual variables like the count, number of outs, score differential, and baserunners. At a high-level **BARTxR** generalizes existing run expectancy measures like **RE24** ([Weinberg, 2014](#)), which computes the average number of runs scored within bins defined by the number of outs and configuration of baserunners. Using a comprehensive cross-validation study, we demonstrate that **BARTxR**’s predictions are more accurate than those of **RE24** and several variants thereof.

The remainder of this paper is organized as follows. We introduce the data and notation used in [Section 2.1](#) before describing our three-step framework for assessing batter decision-making in [Section 2.2](#). Then, we detail our modeling approach in [Section 3](#). In [Section 4](#), we perform a case study of a single batter, Mike Trout, highlighting the types of plate discipline assessments that can be done using our framework. We conclude in [Section 5](#) with a discussion of several extensions of our modeling framework.

2 Data and background

2.1 Data and notation

Our analysis uses pitch-by-pitch tracking data from Major League Baseball’s Statcast database. For each pitch, we observe game-state information (e.g. count, outs, score, baserunners), pitch personnel (e.g. pitcher, batter, fielders, umpire), pitch outcome (e.g. hit, ball, strike), and other game actions (e.g. steal, substitution). Additionally, we observe the horizontal and vertical coordinates of each pitch’s trajectory as it crosses the front edge of home plate. We scraped these data using the **baseballr** package (version 1.2.0; [Petti and Gilani, 2022](#)). We limited our analysis to pitches thrown during regular season games. We fit our strike probability model using all 380,654 pitches that were taken during the 2019 season and we fit our contact probability model using all 341,725 pitches at which batters swung in the 2019 season. We fit our expected runs model (BARTxR) using all 2,853,912 pitches thrown between the 2015–2018. Observe that these three datasets are disjoint.

We use \mathcal{G} to denote information about the state of the game when the pitch was thrown including the count, outs, baserunners, score differential, inning, and whether it is the top or bottom of the inning. We similarly use \mathcal{P} to record the personnel involved in a pitch including the identifies of the batter, catcher, pitcher, and home plate umpire. We additionally include indicators of the batter and pitcher handedness in \mathcal{P} as well as quantitative measures of the batter and pitcher quality, which we describe below. Finally, we denote the location of the pitch (i.e. the `plate_x` and `plate_z` coordinates in Statcast) as its crosses the front edge of home plate with \mathcal{L} .

Let `swing` = 1 if the batter swings and `swing` = 0 if the batter takes the pitch. If the batter decides to swing, let `contact` = 1 if the batter makes contact with the pitch and `contact` = 0 if the batter misses. If the batter does not swing, let `strike` = 1 if the umpire calls a strike and `strike` = 0 if the umpire calls a ball. Let `outcome` = (`contact`, `strike`) be a vector denoting the outcome of the pitch where `contact` = NA if `swing` = 0 and `strike` = NA if `swing` = 1. Let `gstate` denote the game-state category the game moves to after the outcome of the pitch is observed. `gstate` encodes similar information as `outcome` but treats a called strike and a miss as the same game-state.

We quantify batter and pitcher quality using a running estimate of their weighted on based average (wOBA; [Slowinski, 2010b](#)). For each game, we set these quality measures to the be their mean wOBA, averaged over all of their previous games in the season. For batters,

higher wOBA represents higher quality while for pitches, lower wOBA represents higher quality. Note that for the first game, we set each player’s quality measure to be the league average wOBA from the previous season. Our choice to track player quality using a running wOBA estimates follows the example of Brill et al. (2022).

2.2 Determining the optimal decision

To motivate our three-step framework, consider a pitch at the moment just before the batter decides to swing. Suppose that we knew the number of runs the batter’s team would score in the remainder of the half-inning if (a) the batter swings or (b) the batter does not swing. Based on that knowledge, the optimal decision is intuitively the one that leads to scoring more runs. Of course, we are uncertain about these run values at the moment just before the batter swings or takes the pitch. This is because there is uncertainty in the ultimate outcome: if he swings, the batter might make contact or miss and if he takes the pitch, the umpire may call it a ball or a strike. Because of this uncertainty, we propose to determine the optimal swinging decision using the *expected* number of runs the batter’s team will subsequently score, which can be computed by averaging over the uncertainty about these outcomes.

Figure 2 illustrates the four possible outcomes following the batter’s decision. Computing the expected number of runs scored after a swing requires knowledge of (i) the expected number of runs scored after making contact; (ii) the expected number of runs scored after missing; and (iii) the probability of making contact. Similarly, computing the expected number of runs scored after a take requires knowledge of (i) the probability of a called strike; (ii) the expected number of runs scored after a called ball; and (iii) the expected number of runs scored after a called strike.

Formally, let R be the number of runs the batter’s team scores following the pitch. Given the game-state information \mathcal{G} , personnel \mathcal{P} , and location \mathcal{L} of a pitch, we need to compute $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}]$ for both values of `swing` in order to determine the optimal decision for that pitch. Observe that we can decompose the expected runs following a swing or take as

$$\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}] = \sum_{\text{outcome}} \mathbb{P}(\text{outcome}|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing})\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}]. \quad (1)$$

To determine the optimal decision, we introduce `EVdiff` as the difference in expected runs

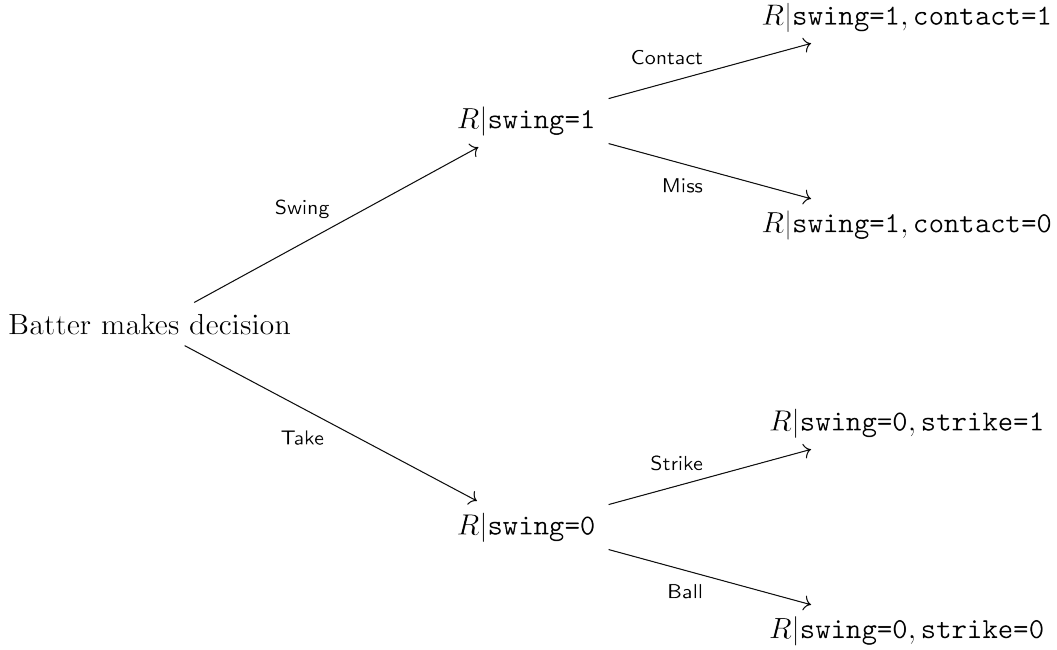


Figure 2: Framework for modeling the outcomes of a pitch.

following a swing and following a take. Formally, we define

$$\text{EVdiff} = \mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing} = 1] - \mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing} = 0]. \quad (2)$$

The sign of EVdiff determines the optimal decision xR_optimal , which is given by

$$\text{xR_optimal} = \begin{cases} \text{swing} = 1, & \text{if } \text{EVdiff} > 0 \\ \text{swing} = 0, & \text{if } \text{EVdiff} \leq 0. \end{cases}$$

Determining the optimal decision for a pitch and subsequently assessing a batter's decision-making requires knowledge of $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}]$ and $\mathbb{P}(\text{outcome}|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing})$. Although we do not know these quantities exactly, we can estimate them using our collected data. We will discuss each estimation problem in Section 3.

2.3 Related work and background on BART

Existing plate discipline metrics. Initial plate discipline metrics (e.g., [Slowinski, 2010a](#)) characterized how often batters swing at pitches thrown outside ($O\text{-Swing}\%$) and inside ($Z\text{-Swing}\%$) the strike zone as well as how often batters made contact with these pitches ($O\text{-Contact}\%$ and $Z\text{-Contact}\%$, respectively). Intuitively, disciplined batters have high $Z\text{-Swing}\%$ and low $O\text{-Swing}\%$. Unfortunately, such metrics implicitly assume that all pitches thrown inside (resp. outside) the strike zone are equally desirable to hit (resp. not hit) which can give a false impression of a batter’s plate discipline. For example, Kevan Smith’s decision in [Figure 1](#) would have a negative impact on his $O\text{-Swing}\%$ even though he had a high probability of making contact. By treating every pitch equally, they fail to account for important contextual factors that impact a batter’s decision to swing such as strike probability (see, e.g., [Arthur, 2014b](#); [Green and Daniels, 2014](#); [Mills, 2014](#)). As a result, they may systematically over-penalize batters who see many pitches near the edges of the strike zone.

To overcome these limitations, recent studies have focused on directly modeling the batter’s actual decision-making process. For instance, [Arthur \(2014a\)](#) first estimates the probability that each pitch is called a strike and then uses those estimates to predict whether a batter will swing at a pitch. They characterize disciplined batters as those whose swing probabilities increase most in response to small increases in strike probability. Unfortunately, this characterization does not consider the downstream impact swing decisions have on metrics like run expectancy or win probability. Towards this end, [Vock and Vock \(2018\)](#) used a causal inference framework to predict how a batter’s batting average, on base percentage, and slugging percentage would change under different counterfactual decision-making strategies. At a high-level, their framework answers the question “What would happen if batter A made swing/take decisions like batter B?” While interesting and informative, the models in [Arthur \(2014a\)](#) and [Vock and Vock \(2018\)](#) make no attempt to determine which decision a batter ought to make. In the context of [Figure 1](#), these models try to predict which path a batter will follow. Our approach, in contrast, attempts to determine which path would most benefit the batter’s team.

Independently of but concurrently to this work, [Mould and Anderson \(2022a,b\)](#) introduced the Expected Additional runs Gained by Looking/swinging Estimate (EAGLE) model to quantify plate discipline. Like us, they also used a tree-structured framework to determine the optimal swing/take decision and fit intermediate strike probability and contact probability models using flexible, non-parametric procedures. We pause here to highlight some

differences between their model and ours. Superficially, our approach differs from theirs in terms of the predictor variables included in these intermediate probability models as well as the specific model fitting procedure used (BART in our case versus XGBoost in their’s). Like us, EAGLE combines predictions from their intermediate probability models with a run expectancy model to compute the expected number of runs a team will score following a swing or take. EAGLE uses a variant of RE24 that accounts for baserunners, count, and outs¹ followed by an *ad hoc* correction for batter quality. Our approach, on the other hand, uses predictions from a much higher-resolution regression-based run expectancy model.

The most substantive difference between our proposal and EAGLE’s lies in uncertainty quantification. Simply put, EAGLE makes no attempt to propagate uncertainties about the intermediately-estimated strike or contact probabilities and run expectancies to their evaluation of batter decision making. In sharp contrast, our Bayesian approach makes such uncertainty quantification and propagation extremely easy (see Section 3.1). We argue that such uncertainty quantification is of paramount importance for evaluating decision-making. Basically, we do not wish to penalize batters for making suboptimal decisions when there is considerable uncertainty about what the optimal decision is.

BART. Initially introduced in the context of nonparametric regression, BART has emerged as an extremely popular “off-the-shelf” modeling tool because it often delivers extremely accurate predictions with reasonably well-calibrated uncertainty estimates without requiring users to (a) pre-specify the parametric form of the regression function and (b) tune any hyperparameters. At a high-level, BART works by approximating unknown functions with sums of binary regression trees and excels at capturing complicated nonlinearities and complex, high-order interaction effects. We believe *a priori* that both strike and contact probabilities are highly non-linear and may depend on complicated interactions between players, umpires, pitch location, and in-game contextual variables like the count or baserunner configuration. Insofar as such nonlinearities and interactions are difficult to specify correctly in a parametric fashion, BART is an especially attractive modeling choice as it does not require us to pre-specify such a parametric form. We fit our three BART models using the **flexBART** package, which permits more flexible modeling with categorical predictors like batter identity that can take on many values. See [Deshpande \(2022\)](#) for more details.

¹We include this model in our cross-validation study in Section 4.1 as RE288

3 Modeling and uncertainty propagation

3.1 Uncertainty propagation

Recall that computing `EVdiff` requires first estimating each term in the summand in Equation (1). By taking a Bayesian approach, we can quantify how uncertainties about these estimates propagate to uncertainty about `EVdiff` and `xR_optimal` in a relatively straightforward fashion. Specifically, because they involve non-overlapping subsets of pitches, we can fit independent Bayesian models to obtain posterior samples of $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}]$ and $\mathbb{P}(\text{outcome}|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing})$ for each outcome. By suitably multiplying and summing these samples according to Equation (1), we obtain posterior samples of $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}]$, from which we can immediately compute posterior samples of `EVdiff` and `xR_optimal`. Given posterior samples of `EVdiff`, we can use the proportion of samples with positive `EVdiff` as an estimated probability that swinging is the optimal decision. We can further quantify how much better (or worse) swinging at the pitch is than taking the pitch using the posterior mean and 90% credible interval (formed using the 5% and 90% sample quantiles) of `EVdiff`.

3.2 Modeling assumptions and fitting

Expected Runs. We need to compute run expectancy following the four pitch outcomes shown in Figure 2. That is, we need to compute $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}]$. A natural approach begins by first partitioning our dataset of 2,853,912 pitches into bins formed by every combination of variables in \mathcal{G} , \mathcal{P} , \mathcal{L} , `swing`, and `outcome` and then computing the average number of runs scored subsequently in each bin. Despite its intuitive appeal, such a binning and averaging procedure is impractical without further simplifying assumptions due to the sheer size of the number of bins. To wit, in 2019 there were a total of 988 batters, 93 umpires, and 113 catchers. Accounting for all combinations of just these three aspects of \mathcal{P} requires over 10 million bins which is greater than the number of pitches in our dataset.

To simplify this estimation, we assume that given the game context, swing decision, and outcome, personnel and pitch location have no predictive effect on expected runs. Formally, we assume that $\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}] = \mathbb{E}[R|\mathcal{G}, \text{swing}, \text{outcome}]$. While this may seem like a rather strong assumption, we note that it is actually weaker than the assumptions underpinning other popular run expectancy models. For example, `RE24` assumes that, given the number of outs and the configuration of baserunners, R is conditionally independent of all other contextual factors, personnel, and pitch location. That is, `RE24` assumes that

$$\mathbb{E}[R|\mathcal{G}, \mathcal{P}, \mathcal{L}, \text{swing}, \text{outcome}] = \mathbb{E}[R|\text{outs}, \text{baserunners}].$$

Under our assumption, we must now compute $\mathbb{E}[R|\mathcal{G}, \text{swing}, \text{outcome}]$. While binning and averaging is now feasible, we instead fit a single BART model to estimate the expected runs following a strike, ball, contact, and miss. Fitting a single model to estimate all four quantities allows us to “borrow strength” from related observations with different outcomes. Such partial pooling may be preferable to separately modeling each outcome since different outcomes can lead to the same game state (e.g. swinging and miss vs. called strike). To model expected runs, we regress R on \mathcal{G} , `swing`, `contact`, `strike`, `outcome`, and `gstate` using pitch data from the 2015 to 2019 MLB seasons.

Event Probabilities. We fit BART models with probit links to estimate the strike probability and contact probabilities as functions of \mathcal{G} , \mathcal{P} , and \mathcal{L} . We fit these models to data from the 2019 MLB season.

4 Results

Code to download and pre-process our data, fit the constituent models, and compute the posterior distributions of `EVdiff` and `xR_optimal` is available at https://github.com/ryanyee3/plate_discipline_code.

4.1 Model validation

We performed several cross-validation studies to understand the predictive accuracy of our BART models for strike probability, contact probability, and run expectancy.

Strike and contact probabilities. We compared our BART models for strike and contact probabilities to two competitors, a model based on generalized additive models (GAMs) and the empirical probabilities of each event using 10-fold cross-validation. We used the `mgcv` package (version 1.8-41; Wood, 2023) to fit GAMs to predict strike and contact probabilities as a smooth function of location with a binomial link. We evaluated each model in terms of mean-squared error (i.e. the Brier score), log-loss, and misclassification error.

Unsurprisingly, both BART and the GAM substantially outperformed the empirical probabilities on both tasks. We further found that BART achieved slightly smaller out-of-sample log-loss than the GAM for strike probabilities (0.16 vs 0.17) and contact probabilities (0.46 vs. 0.47). That the two models achieved similar predictive performance, despite the fact

that our BART models utilized many more predictors than the GAM, suggests that location is, by far, the main driver of strike and contact probabilities. A table with the full results of each fold and each loss metric can be found in Section 6.

Expected runs. We conducted a comprehensive comparison of BARTxR to 14 candidate run expectancy models. We utilized a similar framework to RE24 to fit each candidate model by binning and averaging the number of runs following different game situations (Weinberg, 2014). We consider every combination of count (12 possibilities), outs (3 possibilities), and baserunners (8 possibilities) as predictors in candidate models for a total of seven unique sets of predictors. For each set of predictors we fit two types of models REx and BayesREx where x denotes the number of predictors. REx models are fit in the same way as RE24. BayesREx models are hierarchical Bayesian models of the following form

$$\begin{aligned}
 \bar{\beta} &\sim \mathcal{N}(0, \tau_{\beta}^2) \\
 \tau_{\beta} &\sim \text{half-}t_7 \\
 \sigma^2 &\sim \text{Inverse Gamma} \left(\frac{\nu}{2}, \frac{\nu\lambda}{2} \right) \\
 \beta_{g(i)} &\sim \mathcal{N}(\bar{\beta}, \tau_{\beta}^2) \\
 \tilde{R}_i &\sim \mathcal{N}(\beta_{g(i)}, \sigma^2)
 \end{aligned} \tag{3}$$

where i indexes the pitch, $g(i)$ indexes the bin (i.e., combination of count, out, and/or baserunners), half- t_7 denotes a t -distribution with 7 degrees of freedom truncated to the positive axis, and \tilde{R} is R standardized to have mean zero and variance one. We set ν and λ so that the prior probability on the event $\sigma < 1$ is about 90%. Our decision to specify the model in Equation (3) on the standardized scale and our choices of ν and λ mirror the choices made in default implementations of BART. We complete our prior specification with the weakly informative choice $\tau_{\beta}^2 = 100$. Compared to REx, which estimates run expectancy in each bin independently, BayesREx “borrows strength” across bins.

Figure 3 shows the mean-squared error results of 10-fold cross-validation on each model. We find REx and BayesREx models perform similarly for each combination of predictors. This is perhaps unsurprising given the large number of pitches; basically, the conditional posterior distribution of each $\beta_{g(i)}$ given $\bar{\beta}, \sigma^2$, and τ_{β} is sharply concentrated around the average values of the \tilde{R}_i ’s in the bin. We find that the models which account for outs and baserunners perform best. BARTxR outperforms all other models in both training and testing samples (average RMSE, BARTxR: 0.95, RE288: 0.96, BayesRE288: 0.96).

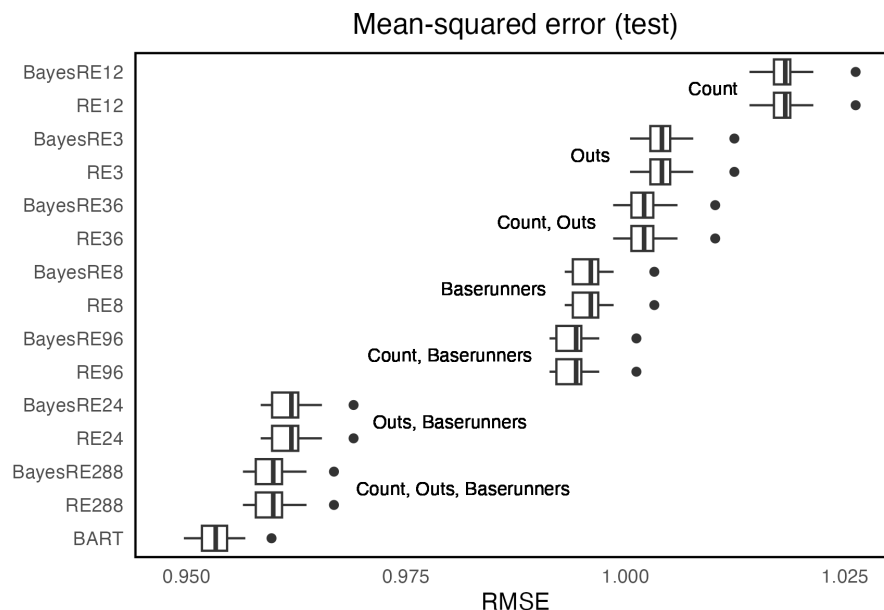


Figure 3: Out-of-sample root mean-squared error for BARTxR, seven REx models, and seven BayesREx models across 10 training/testing splits.

4.2 Batter evaluation case study

We now illustrate the type of plate discipline analysis facilitated by our modeling efforts with a case study about Mike Trout in the 2019 MLB season. For each pitch Trout faced in 2019, we generated samples of `EVdiff` and determined the `xR_optimal` decision based on the posterior mean of `EVdiff`. We quantified our uncertainty by computing the `xR_optimal` decision for each individual sample; that is, we computed the proportion of samples in which the `xR_optimal` decision is to swing. Figure 4 shows the results of these calculations for every pitch faced by Trout in the 2019 broken down by the actual decision and the `xR_optimal` decision. In Figure 4a, the darker the shading, the greater the difference in the posterior means of expected runs from swinging and the expected runs from taking. In Figure 4b, the darker the shading, the more posterior certainty we have in the `xR_optimal` decision.

In Figure 4a and Figure 4b, the large number of pitches in panels (a) and (d) reveal that Trout's actual decision very often matched our model's `xR_optimal` decision. Moreover, these decision tended to agree with the conventional wisdom that batters should swing at pitches inside the strike zone and take pitches outside the strike zone. In panels (a) and (d) of the figures, we find, for instance, that Trout tended not to swing at pitches thrown outside the strike zone and rarely took pitches inside the zone. The dark shading of pitches

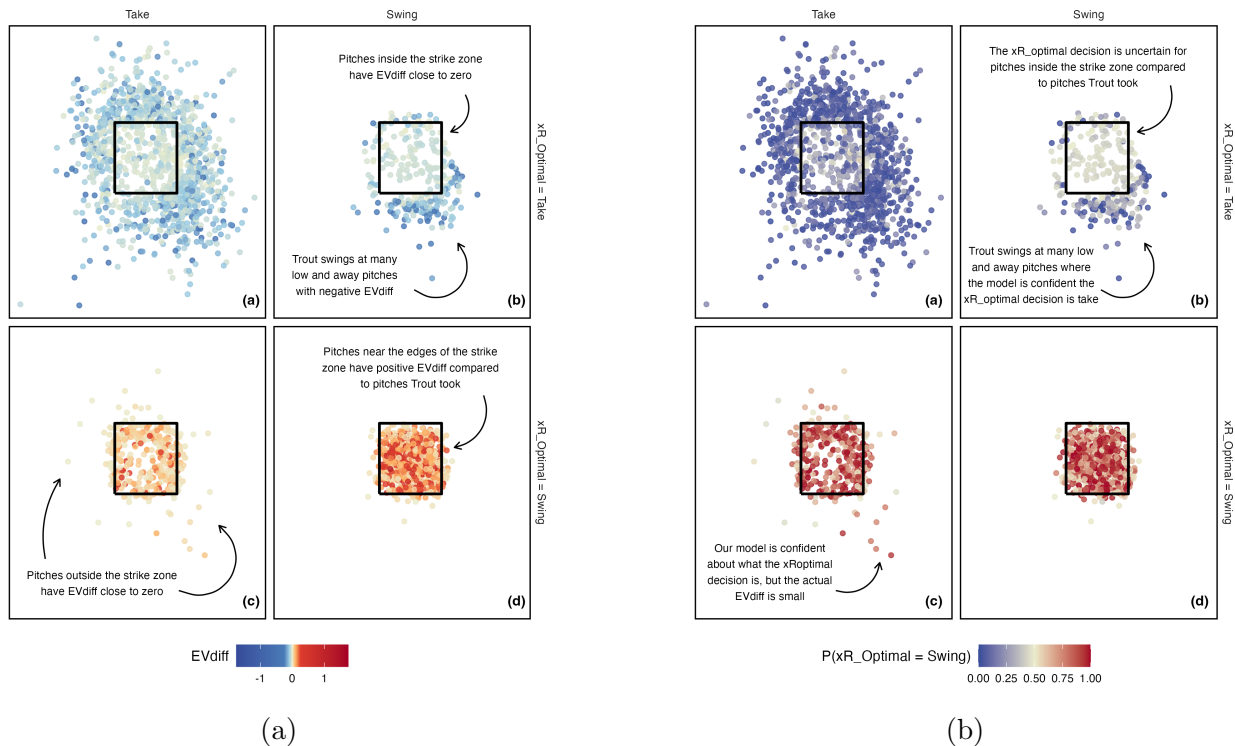


Figure 4: Pitches faced by Mike Trout in the 2019 MLB season where he took (left) and swung (right) and the $xR_optimal$ decision is to take (top) and swing (bottom). Pitches in (a) are shaded based on $EVdiff$ and pitches in (b) are shaded based on $\mathbb{P}(xR_optimal = swing)$. The darker the shading, the larger the $EVdiff$ (a) and our certainty that the $xR_optimal$ decision is to swing (b). All plots are from the perspective of the home plate umpire.

outside the strike zone in panel (a) illustrate our model’s high degree of certainty that the $xR_optimal$ decision is to take the pitch. Similarly, the dark shading of pitches inside the strike zone in panel (d) illustrate our model’s certainty that the $xR_optimal$ decision is to swing.

Perhaps more interesting are those pitches where our model’s $xR_optimal$ decision deviated from the conventional wisdom and Trout’s actual decision-making (panels (b) and (c) in Figure 4a and Figure 4b). In panel (b), for instance, we see that Trout swung at many pitches inside the strike zone for which our model determined the $xR_optimal$ decision was to take. As suggested by the relatively light shading of these pitches, we found on further inspection that our model was very uncertain about the $xR_optimal$ decision. In fact, the posterior distributions of $EVdiff$ for these pitches tended to be nearly symmetric and tightly concentrated around zero. Because of the uncertainty in the $xR_optimal$ decision and rel-

atively small magnitude of EV_{diff} for these pitches, we would not classify Trout’s decision to swing at these pitches as bad decisions per se.

In panel (c), however, we find that Trout took several pitches that passed inside the strike zone for which our model determined, with relatively high certainty (as evidenced by the dark shading), that the `xR_optimal` decision was to swing. For these pitches, the posterior distributions of EV_{diff} were largely concentrated on positive values. By taking these pitches, our model suggests that Trout cost his team in terms of expected runs.

Of additional interest are the low-and-away pitches in panel (c) of Figure 4b where our model is very confident that the `xR_optimal` decision is to swing. We found that the posterior mean EV_{diff} for these pitches is close to zero, suggesting that the cost of making a sub-optimal decision on these pitches is very small. Nevertheless, it is interesting that our model is so confident that the `xR_optimal` decision is to swing on these pitches thrown well outside the strike zone, when it is similarly confident that the `xR_optimal` decision is to take pitches thrown in similar locations in panels (a) and (b) of Figure 4b. We speculate that the difference is due, at least in part, to differences in the game contexts in which these pitches were thrown.

To probe this possibility, we visualized all pitches faced by Trout broken down by the number of outs and the number of baserunners in Figure 5. Pitches are colored in Figure 5 based on the posterior probability that the `xR_optimal` decision is to swing. We see that the low-and-away pitches where the `xR_optimal` decision is to swing (pitches of interest) occur when there are two outs and few baserunners while the low-and-away pitches where the `xR_optimal` decision is to take occur when there are no outs. Such a finding is, we argue, intuitive, because the expected runs from a positive outcome will be similar to the expected runs from a negative outcome when there are two outs and few baserunners.

Finally, while we have focused on Trout in this paper, we can conduct such analyses for any batter using our model results. We have created an interactive Shiny app (version 1.7.2; Chang et al., 2022) that performs such visual analysis for all batters who faced at least 1,000 pitches in the 2019 MLB season. The app is available at https://ryanyee3.shinyapps.io/batter_evaluation_app/.

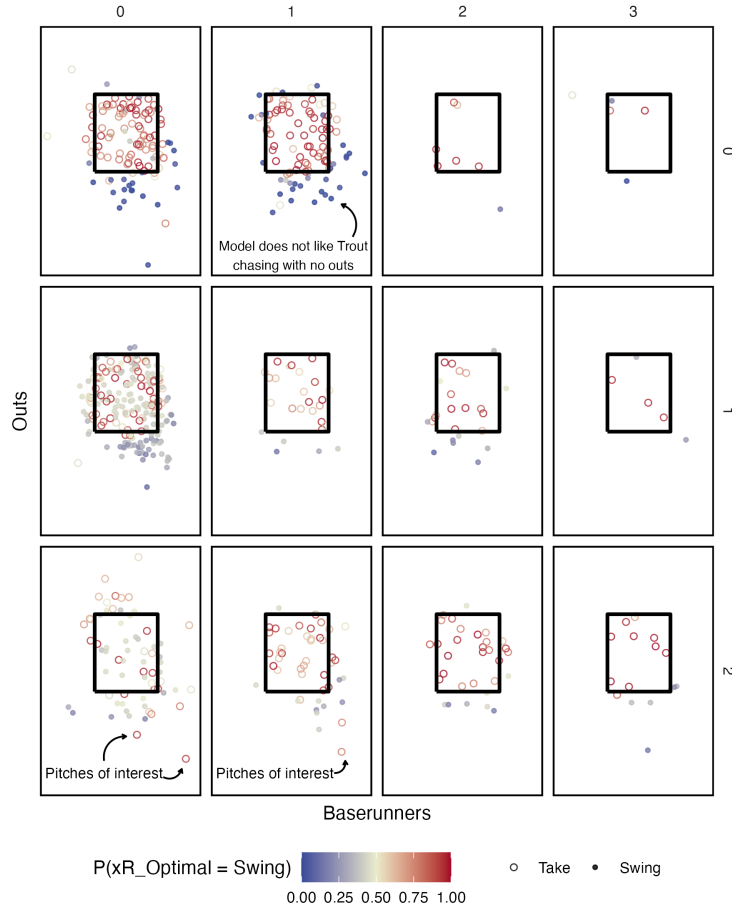


Figure 5: Pitches faced by Mike Trout in the 2019 MLB season where his swing decision conflicts with the `xR_optimal` decision. Pitches are shaded based on the proportion of samples where the expected runs from swinging is greater than the expected runs from taking.

4.3 Summary metrics

We can complement our visual analysis using several aggregate metrics. For instance, we can compute the proportion of pitches where the batter actually makes the `xR_optimal` decision. To account for our uncertainty in the `xR_optimal` decision, we can calculate this proportion for every posterior sample of `xR_optimal`, and compute a credible interval for this metric. We find that the average batter made the `xR_optimal` decision 69.4% of the time in the 2019 MLB season. For reference, based on the traditional heuristic that a disciplined swing decision is to swing at pitches inside the strike zone and take pitches outside the strike zone, batters made the disciplined decision 68.4% of the time in the 2019 MLB season. We found

that Jonathan Luplow had the highest proportion of `xR_optimal` swing decisions (76.8%, 90% credible interval: [46.1%, 88.7%]) while Jorge Alfaro had the lowest proportion of `xR_optimal` swing decisions (61.4%, 90% credible interval: [37.7%, 83.2%]). These numbers closely align with traditional plate discipline metrics that credit Luplow and Alfaro with making the disciplined decision 76.4% and 60.7% of the time, respectively.

We can also measure the total impact a batter’s decisions have on expected runs over an entire season by taking the difference in expected runs between the decision the batter made and the alternative decision and summing over every pitch a batter faces. In terms of the plots in Figure 4, let $\sum_p \text{EVdiff}$ be the sum of `EVdiff` of all the pitches shown in panel p . Then the expected runs added over an entire season could be computed as

$$\left(\sum_d \text{EVdiff} - \sum_c \text{EVdiff}\right) - \left(\sum_a \text{EVdiff} - \sum_b \text{EVdiff}\right).$$

Figure 6 shows a histogram of the posterior mean of expected runs added for all batters who faced at least 1000 pitches in the 2019 season. We find the difference between the best and worst batters according to this metric is not large — only about 0.1 expected runs. We further find that there is considerable overlap in the 90% credible intervals of the expected runs added (see Figure 8). For instance, we estimate the top batter, Jordan Luplow, on average adds 0.08 expected runs per pitch due to his decision-making with a 90% credible interval of [-0.06, 0.21] while the worst batter, Javier Baez, adds 0.03 expected runs per pitch on average with a 90% credible interval of [-0.13, 0.18].

While informative, the expected runs added metric does face some limitations. We argue that if a batter faces many low-leverage pitches (i.e. pitches where $|\text{EVdiff}|$ is close to zero) will have fewer opportunities to pick up added runs than batter’s facing many high-leverage pitches (i.e. pitches where $|\text{EVdiff}|$ is large). To account for this, we can calculate runs lost as the minimum of 0 and the runs added from a pitch. To motivate this metric, we argue that when a batter faces a low-leverage pitch, it does not matter what decision he makes since it may be unclear what the best decision is. For example, in Figure 4 panel (b) there are many pitches inside the strike zone that Trout swings at that, while the `xR_optimal` decision is to take, there is still considerable uncertainty in the `xR_optimal` decision. Since the most a batter can be punished on a given pitch is the expected runs lost by not making the `xR_optimal` decision, on low-leverage pitches the batter will either not be penalized or be penalized by a small amount. When a batter faces a high-leverage pitch, the `xR_optimal`

decision should be obvious; and, under our framework, we expect a player with good plate discipline to make the optimal decision, so a batter that does not make the `xR_optimal` decision on these pitches will be severely punished. In terms of the plots in Figure 4, runs lost is calculated as

$$\sum_c \text{EVdiff} - \sum_b \text{EVdiff}.$$

Results of expected runs lost are similar to the results of expected runs added: the absolute differences between the best and worst batters are small (see Figure 6) with considerable overlap in credible intervals (see ??). We estimate the top batter, Andrew McCutchen, has an average loss of 0.03 expected runs per pitch with a 90% credible interval of [0.01, 0.09] while the worst batter, Jorge Alfaro, lost 0.07 expected runs per pitch on average with a 90% credible interval of [0.01, 0.20].

For all three of these metrics, we find that there are a few batters that do much better than everyone else and a few that much worse and the majority are somewhere in-between. Figure 6 shows the distribution of these metrics for batters that faced at least 1000 pitches in the 2019 MLB season.

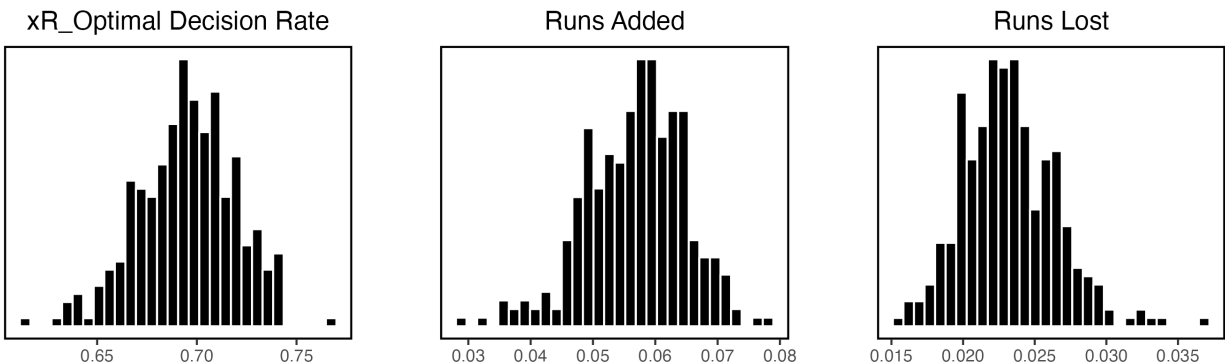


Figure 6: Distributions of posterior means of batters who faced at least 1000 pitches in the 2019 MLB season for each summary metric.

While it is tempting to use these metrics to compare the plate discipline of different batters, note that these metrics are confounded by the location of pitches a batter faces, so they are not necessarily directly comparable across batters. We will return to this point in Section 5.

5 Discussion

We developed a three-step framework for estimating the optimal swing/take decision of batters in Major League Baseball. In the first step of our framework, we estimate, for a given pitch, (i) the probability the umpire will call a strike; (ii) the probability a batter will make contact; and (iii) the expected number of runs the batting team will score in the remainder of the inning after a ball, strike, contact, and miss. Then, we combine these estimates to calculate the expected runs following a swing and a take decision. Finally, we determine the `xR_optimal` decision to be the one that leads to more expected runs. We adopt a Bayesian approach which allows us to propagate uncertainty from each intermediate estimate to our evaluation of the `xR_optimal` decision.

Our findings have several implications for Major League Baseball teams. First, we can determine those batters who consistently make `xR_optimal` swing decisions, which can aid in player evaluation. We can also identify pitch locations where a batter consistently makes suboptimal swing decisions. Batters can use this information to make adjustments to their decision-making, while pitchers can use this information to target locations where a batter is likely to make a decision that costs his team runs. Finally, we can identify locations where pitches are likely to lead to a significant increase in expected runs (i.e. locations with high `EVdiff`). Pitchers can use this information to avoid these locations to minimize the opportunities for the batting team to score.

We found that across the leagues, batters made the `xR_optimal` decision 69.4% of the time according to our model, which aligns with the findings of other plate discipline studies. Critically, we should ask ourselves why are our results so similar to the simpler traditional plate discipline metrics. One explanation is that, despite their somewhat heavy-handed assumptions (e.g. one should swing at every pitch in the strike zone and take every pitch outside the zone), traditional metrics offer a reasonably accurate approximation of our more fine-grained model.

Like umpires, batters are not robots: just as umpires do not strictly adhere to the rulebook definition of the strike zone, batters almost certainly consider more than just expected runs when making swing decisions. The fact that batters deviated from the `xR_optimal` decision about 30% of the time suggests that batters do not always consider expected runs when making swing decisions. For example, a player chasing a home run record (e.g., Aaron Judge at the end of the 2022 season) may make swing decisions to maximize expected home runs as

opposed to expected runs. While these decisions are not `xR_optimal`, they might be optimal when viewed through the lens of expected home runs.

5.1 Limitations and future work

We presented an analysis of plate discipline that asserts the “disciplined” swing decision is the one that leads to a greater number of expected runs. While we believe our assertion is reasonable, an honest argument could be made that this is an incorrect characterization of a “disciplined” swing decision. In such cases, our modular framework gives us the flexibility to use any other objective to evaluate batters. For example, we could replace R in each branch of Figure 2 with another team (e.g., win probability) or individual (e.g. home runs, wOBA, OPS) outcome and following the same strategy: estimate intermediate probability and conditional expectations, compute the expected outcome following a swing and following a take, and determine the optimal decision.

While we used BART to fit the three models in the first step of our framework, other choices are possible. Indeed, we selected BART for its ease-of-use: we did not have to manually pre-specify the functional forms of the called strike and contact probabilities and run expectancy, which we suspected depend on complicated non-linearities and interactions. Although our BART-based models slightly outperformed parametric alternatives, these differences in out-of-sample predictive power were not especially large. For instance, we found that our BART model of called strike and contact probabilities, which accounted for location, game state, and pitch personnel, achieved very similar predictive performance as a generalized additive model that only accounted for location. Such a finding, to us, indicates that pitch location is, by far, the main driver of strike and contact probabilities, with other predictors like outs or count or player identifies contributing little additional predictive power. We also found that the difference in out-of-sample RMSE between our BART-based run expectancy model, `BARTxR`, and the best-performing alternative based on binning and averaging was about 0.01 runs. While this is a small difference on a per-pitch basis, these differences can magnify over the course of an entire season.

More substantively, we only consider two possible outcomes following a swing decision, a miss or contact. We could extend our framework to account for more post-swing outcomes like miss, foul, out, single, double, triple, and home run in one of two ways. First, we could replace our binary contact probability model with Murray (2021)’s multinomial logistic BART model. In the context of Figure 2, this would involve replacing the `contact` branch

with several branches. Alternatively, we could augment our existing framework with a further model these outcomes conditional on making contact. That is, we could add additional child branches to the `contact` branch of Figure 2, one for each potential outcome following contact. We note that the developers of EAGLE pursued the first strategy. We suspect, however, that the second approach would lead to less uncertainty about the `xR_optimal` decision as the overall accuracy of our contact/miss model is better than the reported accuracy of EAGLE’s multi-outcome model.

Beyond a pitch-by-pitch visual assessment, we introduced three aggregate metrics that tracked the proportion of times batters made the `xR_optimal` decision and the expected number of runs added or lost due to plate discipline across an entire season. Although it is tempting to compare these metrics across players, these metrics are confounded by pitch location and the context in which players see different pitches. Simply put, because the distribution of pitches one batter sees may differ from the distribution seen by another batter, it is difficult to directly compare their aggregated plate discipline metrics. To overcome such confounding, it is tempting to first marginalize `EVdiff` over \mathcal{L} and \mathcal{G} before computing in a manner similar to Jensen et al. (2009)’s spatially aggregate fielding evaluation. While such spatially and contextually aggregated plate discipline metrics are intuitively appealing, computing them introduces considerable computational challenges. Basically, to marginalize over \mathcal{L} and \mathcal{G} , we would have to make predictions for every combination of batter and pitch in our dataset. We leave efficient computation of such marginal measures to future work.

Relatedly, our strike and contact probability models both condition on the location of the pitch as it crosses the front edge of home plate. Insofar as batters decide to swing or take the pitch before it reaches home plate, one could reasonably argue that such location information is unavailable to the batter. Rather than omit pitch location from our strike and contact probability models, it would be interesting to fit a more realistic model that conditions on the location of the pitch at a point between the pitcher’s mound and home plate.

References

- Arthur, R. (2014a). Moonshot: The new best way to measure plate discipline. <https://www.baseballprospectus.com/news/article/25008/moonshot-the-new-best-way-to-measure-plate-discipline/>.
- Arthur, R. (2014b). Moonshot: The victims of a bad strike zone. <https://www.baseballprospectus.com/news/article/24862/moonshot-the-victims-of-a-bad-strike-zone/>.

- Brill, R. S., Deshpande, S. K., and Wyner, A. J. (2022). A bayesian analysis of the time through the order penalty in baseball. *arXiv preprint arXiv:2210.06724*.
- Chang, W., Cheng, J., Allaire, J., Sievert, C., Schloerke, B., Xie, Y., Allen, J., McPherson, J., Dipert, A., and Borges, B. (2022). *shiny: Web application framework for R*.
- Chen, D. L., Moskowitz, T. J., and Shue, K. (2016). Decision mkaing under the Gambler’s Fallacy: Evidence from asylum judges, loan officers, and baseball umpires. *The Quarterly Journal of Economics*, 131(3):1181–1242.
- Chipman, H. A., George, E. I., and McCulloch, R. E. (2010). BART: Bayesian additive regression trees. *Annals of Applied Statistics*, 4(1):266–298.
- Deshpande, S. K. (2022). A new BART prior for flexible modeling with categorical predictors. *arXiv preprint arXiv:2211.04459*.
- Deshpande, S. K. and Wyner, A. J. (2017). A hierarchical Bayesian model of pitch framing. *Journal of Quantitative Analysis in Sports*, 13(3):95–112.
- Green, E. and Daniels, D. (2022). Bayesian instinct. *SSRN preprint 2916929*.
- Green, E. and Daniels, D. P. (2014). What does it take to call a strike? three biases in umpire decision making. In *MIT Sloan Sports Analytics Conference*.
- Jensen, S. T., Shirley, K. E., and Wyner, A. J. (2009). Bayesball: A bayesian hierarchical model for evaluating fielding in major league baseball. *Annals of Applied Statistics*, 3(2):491–520.
- Kim, J. W. and King, B. G. (2014). Seeing stars: Matthew effects and status bias in Major League Baseball umpiring. *Management Science*, 60(11):2619–2644.
- Lindbergh, B. (2013). The art of pitch framing. <http://grantland.com/features/studying-art-pitch-framing-catchers-such-francisco-cervelli-chris-stewart-jose-molina-others/>.
- Marchi, M. (2011). Evaluating catchers: Quantifying the framing pitches skill. <https://tft.fangraphs.com/evaluating-catchers-quantifying-the-framing-pitches-skill/>.
- Mills, B. M. (2014). Social pressure at the plate: Inequality aversion, status, and mere exposure. *Managerial and Decision Economics*, 35(6):387–403.
- Mills, B. M. (2017). Technological innovations in monitoring and evaluation: evidence of performance impacts among Major League Baseball umpires. *Labour Economics*, 46:189–199.
- Mould, J. and Anderson, D. (2022a). Quantifying hitter plate discipline with eagle: Part 1. <https://www.baseballprospectus.com/news/article/74173/quantifying-hitter-plate-discipline-with-eagle-part-1/>.
- Mould, J. and Anderson, D. (2022b). Quantifying hitter plate discipline with eagle: Part 2. <https://www.baseballprospectus.com/news/article/74214/quantifying-hitter-plate-discipline-with-eagle-part-2/>.

- Murray, J. S. (2021). Log-linear Bayesian additive regression trees for multinomial logistic and count regression. *Journal of the American Statistical Association*, 116(534):756–769.
- Petti, B. and Gilani, S. (2022). *baseballr: Acquiring and analyzing baseball data*.
- Slowinski, P. (2010a). Plate Discipline. <https://library.fangraphs.com/offense/plate-discipline/>.
- Slowinski, P. (2010b). wOBA. <https://library.fangraphs.com/offense/woba/>.
- Vock, D. M. and Vock, L. F. B. (2018). Estimating the effect of plate discipline using a causal inference framework: an application of the G-computation algorithm. *Journal of Quantitative Analysis in Sports*, 14(2):37–66.
- Weinberg, N. (2014). Re24. <https://library.fangraphs.com/misc/re24/>.
- Wood, S. (2023). *mgcv: Mixed GAM Computation Vehicle with Automatic Smoothness Estimation*.

6 Appendix

6.1 Cross-validation results

Table 1 reports the out-of-sample mean square error (i.e. Brier score), log-loss, and misclassification rate for our BART model, a generalized additive model that only accounts for pitch location, and the empirical probabilities for each fold in our cross-validation experiment. Table 2 is analogous but for the contact probabilities. Observe that our BART models outperforms the GAMs on every fold.

Table 1: Out-of-sample predictive errors for computing strike probabilities on each cross-validation fold.

Fold	MSE			Log-loss			Misclassification rate		
	BART	GAM	EP	BART	GAM	EP	BART	GAM	EP
1	0.050	0.052	0.214	0.163	0.168	0.619	0.071	0.074	0.310
2	0.050	0.052	0.214	0.161	0.166	0.620	0.069	0.073	0.311
3	0.049	0.051	0.216	0.159	0.165	0.622	0.069	0.072	0.314
4	0.050	0.052	0.216	0.164	0.168	0.624	0.072	0.073	0.316
5	0.051	0.053	0.215	0.163	0.168	0.622	0.072	0.074	0.314
6	0.050	0.052	0.214	0.162	0.167	0.619	0.070	0.074	0.310
7	0.050	0.052	0.216	0.163	0.166	0.623	0.071	0.074	0.315
8	0.050	0.052	0.213	0.161	0.166	0.617	0.070	0.073	0.308
9	0.050	0.053	0.215	0.163	0.168	0.622	0.070	0.074	0.312
10	0.050	0.052	0.214	0.162	0.165	0.619	0.070	0.073	0.310

Table 2: Out-of-sample predictive errors for computing contact probabilities on each cross-validation fold.

Fold	RMSE			Log-loss			Misclassification rate		
	BART	GAM	EP	BART	GAM	EP	BART	GAM	EP
1	0.149	0.152	0.182	0.467	0.476	0.549	0.200	0.210	0.239
2	0.149	0.151	0.182	0.466	0.474	0.550	0.201	0.204	0.239
3	0.146	0.149	0.180	0.460	0.468	0.545	0.197	0.200	0.235
4	0.147	0.151	0.181	0.462	0.472	0.547	0.199	0.202	0.236
5	0.147	0.150	0.181	0.463	0.471	0.548	0.198	0.202	0.238
6	0.150	0.154	0.184	0.470	0.481	0.554	0.204	0.210	0.242
7	0.149	0.152	0.182	0.468	0.476	0.550	0.202	0.204	0.239
8	0.149	0.151	0.182	0.467	0.474	0.549	0.201	0.204	0.239
9	0.147	0.151	0.180	0.463	0.473	0.546	0.199	0.202	0.236
10	0.147	0.150	0.181	0.463	0.471	0.547	0.200	0.203	0.237

6.2 Additional figures

Figure 7 shows boxplots of the proportion of pitches where batter’s made the `xR_optimal` decision for each posterior sample for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season. There is considerable uncertainty in the proportion of pitches where batter’s made the `xR_optimal` decision, and the overlap between batters demonstrates the difficulty in differentiating batters.

Figure 8 shows boxplots of the runs added metric for each posterior sample for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season. The boxplots show considerable uncertainty in these estimates, and overlap in distributions between batters demonstrates the difficulty in differentiating batters based on this metric.

Figure 9 shows boxplots of the runs lost metric for each posterior sample for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season. There is considerable uncertainty in these estimates, and overlap in distributions between batters demonstrates the difficulty in differentiating batters based on this metric.

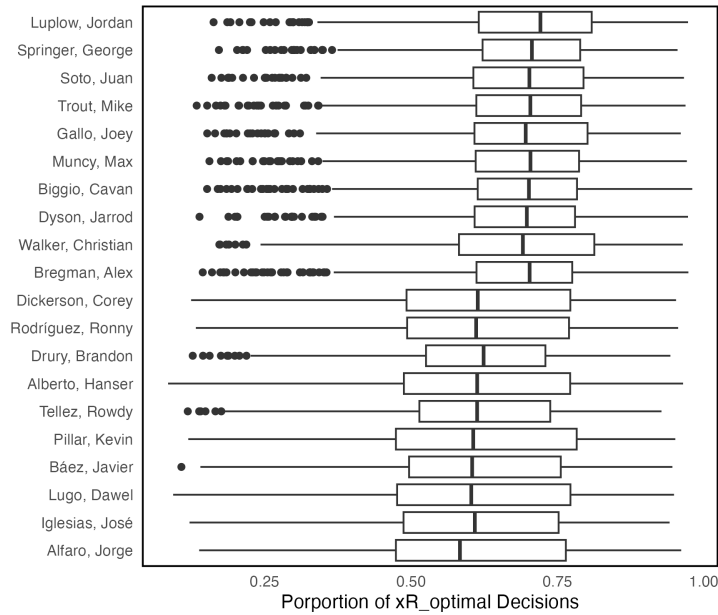


Figure 7: Boxplots of posterior samples of the proportion of pitches where the batter made the `xR_optimal` decision for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season.

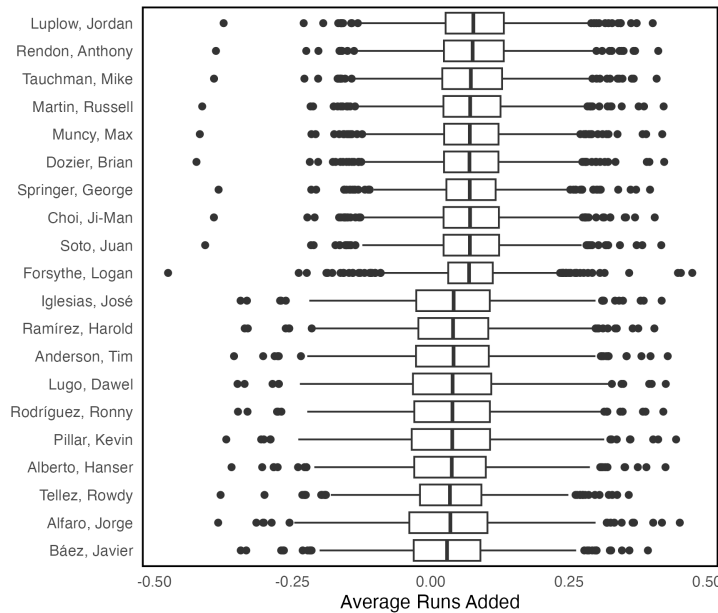


Figure 8: Boxplots of posterior samples of the runs added metric for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season.

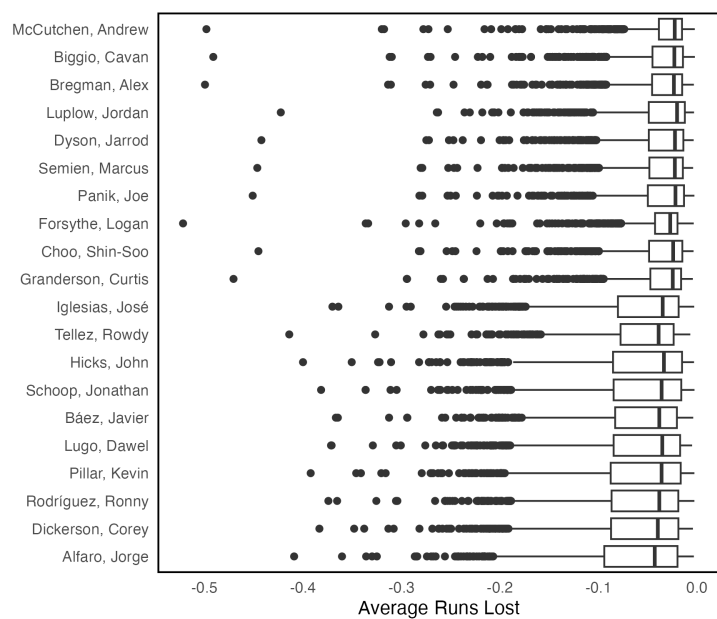


Figure 9: Boxplots of posterior samples of the runs lost metric for the top 10 and bottom 10 batters who faced at least 1000 pitches in the 2019 MLB regular season.