

# Example of the Research Process:

## Rethinking WAR for Starting Pitchers

Suppose I want to start a sports analytics research project.

But, I don't have any ideas right now.

A fantastic way to get started is to simply read about something you're interested about.

Perhaps you were listening to a podcast on which someone mentioned that Roger Clemens has the most career Wins Above Replacement (WAR) of all time, 133.7, according to FanGraphs.

You may have also heard that Pedro Martinez in 1999 has the highest single season WAR of all time, 11.6, according to FanGraphs.

You may think that WAR is a really cool concept, and it makes some intuitive sense

over why it seems like a nice way to evaluate pitchers, and more generally, all players.

Wins Above Replacement — Replace a player with a replacement-level player (e.g. the best guy you could get on waivers), how many fewer wins would the team have, assuming average teammates and opponents?

Historical WAR — how many wins above replacement did a player have last season?

Predictive WAR — how many wins above replacement will a player have next season?

These 2 ideas are often conflated, but they should be separated.

Implementation take a player's observed performance, ignoring/adjusting for things that he is not responsible for, and map that to Wins

Say you don't know the math behind WAR, although you are curious to learn. So, you **Read**.

The most widely used/accepted public WAR implementations are from FanGraphs and Baseball Reference.

FanGraphs WAR for pitchers:

<https://library.fangraphs.com/war/calculating-war-pitchers/>

Baseball Reference WAR for pitchers:

[https://www.baseball-reference.com/about/war\\_explained\\_pitch.shtml](https://www.baseball-reference.com/about/war_explained_pitch.shtml)

When you Read about WAR for pitchers, a few things catch your eye:

\* WAR involves mapping a pitcher's performance  
(e.g., FIP for FanGraphs  
xRA for Baseball Reference)  
to Wins

$$\text{ifFIP} = \frac{13 \cdot \text{HR} + 3 \cdot (\text{BB} + \text{HBP}) - 2 \cdot (\text{K} + \text{IFFB})}{\text{IP}} + C$$

### Fielding Independent Pitching (with Infield Flies!)

The first thing you need to do to calculate a pitcher's WAR is to calculate their FIP. Unfortunately for those of you playing along at home, you can't simply take the pitcher's FIP from their player page because we treat **infield fly balls (IFFB) as strikeouts for the purposes of WAR** but not for the general FIP calculation found on the player's page. We'll call this ifFIP to avoid confusion. Here is the formula:

$$\text{ifFIP} = ((13 * \text{HR}) + (3 * (\text{BB} + \text{HBP})) - (2 * (\text{K} + \text{IFFB}))) / \text{IP} + \text{ifFIP constant}$$

This is the traditional FIP formula, but with IFFB added in as strikeouts. However, keep in mind that you also need to calculate a special ifFIP constant and can't just grab "cFIP" from our guts page.

$$\text{ifFIP Constant} = \lg \text{ERA} - (((13 * \lg \text{HR}) + (3 * (\lg \text{BB} + \lg \text{HBP})) - (2 * (\lg \text{K} + \lg \text{IFFB}))) / \lg \text{IP})$$

$xRA$  = expected runs allowed  
= ignoring the ordering e.g. 1B, out, out, 1B, HR, out  
v.s. HR, 1B, 1B, out, out, out

and just using the events 1 HR, 2 1B, 3 out,  
what is the expected runs allowed of the inning?

- \* there are a series of convoluted adjustments on top of the base metric
  - (e.g., league adjustment  
team defense adjustment)
- \* WAR involves mapping a pitcher's performance averaged over the entire season into wins
  - if FIP: divides by IP
  - $xRA$ : cumulative seasonal  $xRA$

Thoughts: averaging pitcher performance over the course of a season  
Seems weird

Let's explore some implications of this modeling assumption.

Ex

game	1	2	3	4	5	6	total
earned runs	0	10	1	2	1	1	15
innings pitched	9	4	6	7	8	7	41

Table 1: Max Scherzer's performance over six games prior to the 2014 All Star break.

4 dominant performances  $\rightarrow \geq 4$  wins

$$\frac{15 \text{ Runs}}{41 \text{ IP}} \times 9 \frac{\text{innings}}{\text{game}} = \frac{3.66 \text{ Runs}}{\text{Complete game}}$$

$3.66 \frac{\text{Runs}}{\text{Complete game}} \approx 0.55 \frac{\text{Win probability}}{\text{Complete game}}$

$\rightarrow \approx 3.30$  wins

over 6 games

Big difference b/t  $\geq 4$  and 3.3 wins !

Ex Would you rather have pitcher A or pitcher B?

A: 5 Runs in each game

B: Alternates b/t 10 and 0 Runs  
in each complete game

All else the same, existing WAR methodologies value these 2 pitchers the exact same.

Would rather have pitcher B though...

Ex Pitcher A: alternates between allowing 7 and 0  
Runs per complete game

Pitcher B: alternates between allowing 14 and 0  
Runs per complete game

Existing WAR:  $A \sim 3.5$  runs/game  
 $B \sim 7$  runs/game

$$A > B$$

"Real" WAR: Both A and B win half of  
their games.

$$A \approx B$$

"You can only lose a game once"

"Not all Runs have the same value"

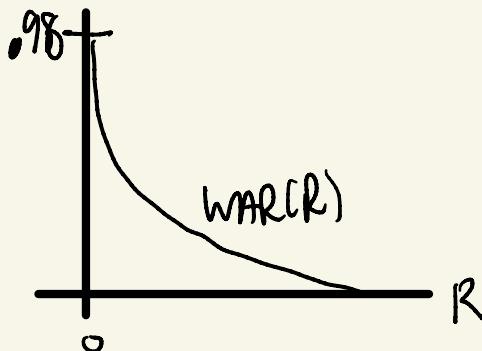
→ the 8<sup>th</sup> Run allowed in a game  
is "worth" less than the 1<sup>st</sup>

→ the marginal difference in Win  
Probability between allowing the 7<sup>th</sup> vs. 8<sup>th</sup>  
Run is less than the marginal  
difference in game WAR between allowing

the 1<sup>st</sup> and 2<sup>nd</sup> Run,  
since you're essentially already  
lost the game

→ if  $R \mapsto \text{WAR}(R)$  is game WAR  
as a function of Runs allowed,  
then  $\text{WAR}(R) - \text{WAR}(R-1)$   
gets smaller as  $R$  gets bigger

→  $R \mapsto \text{WAR}(R)$  should be convex



Problem : Averaging pitcher performance over the course of the season is clearly wrong

- it ignores the game-by-game variance in pitcher performance
- $R \mapsto \text{WAR}(R)$  should be convex  
(i.e., not all Runs should have the same value)
- a win is the fundamental result of a game,  
not a season

**Goal:** Fix this problem.

Make one incremental improvement [Research].

→ calculate historical WAR in each individual game  
seasonal) WAR is the sum of game WAR.

Task Game WAR for starting pitchers

How to do this?

English → Math

(TM)

{  
WAR    WAR = wins above Replacement  
Wins     $W =$  How many wins did Scherzer  
contribute in this game?  
Above Replacement  $W_{Rep} =$  How many wins a  
replacement-level pitcher  
would've contributed  
 $WAR = W - W_{Rep}$

One step at a time

Begin with Wins  $W$

Wins How many wins did Scherzer contribute in this game?

Math: win probability

Pitcher valuation: we only want to judge Scherzer using things he's responsible for

Scherzer's game performance:	Runs Allowed	R	because winning a game is defined by Runs
	exit inning	I	
	exit base-state	S	
	exit outs	O	

Confounders, e.g.

Variables that affect his performance:

PARK

opposing team's batting quality

his team's fielding quality

contextual

variables that affect the win probability

league (NL vs. AL), season

Variables that don't affect his performance

and so we shouldn't judge him with:

his team's batters/opposing team's defense

## Start Simple

Begin with the easiest version of the task.  
Then, iterate on top of that.

\* Begin just with Scherzer's observed performance.  
Adjust for confounders later.

Task given Scherzer's performance,  
what's his team's win probability  
when he exits the game?

Context-neutral: assume league-average offenses, defenses,  
ignore his own team's runs scored

Start simple: assume he finishes the  
inning, so ignore (S, O)

Model the function

$$f = f(I, R) =$$

assuming both teams have  
league-average offenses,  
compute the probability a team  
wins a game after giving up R  
runs through I complete innings

Since  $f(I, R)$  can be visualized as a 2D grid,  
we name our WAR Grid **WAR**.

This is the simplest version of the question, and it is nontrivial.

# Machine Learning vs. Mathematical Models

- Machine Learning / Statistical models are fit from historical data
- Mathematical Models are equations written on paper

## Models fit from historical data

- What data do we need?
- get data
- empirical grid
- logistic regression — No
- XGBoost grid

## Mathematical Models

- Poisson( $\lambda$ ) model
- Poisson( $\lambda_x$ ), Poisson( $\lambda_y$ ) model
- overdispersion hyperparameter K

See Code

## A Estimating $f$ using a mathematical, not a statistical, model

In this Section, we detail our modeling process for estimating the grid function  $f = f(I, R)$  which, assuming both teams have randomly drawn offenses, computes the probability a team wins a game after giving up  $R$  runs through  $I$  complete innings. In particular, we compare statistical models fit from observational data to mathematical probability models, which are superior.

To account for different run environments across different seasons and leagues (NL vs. AL), we estimate a different grid for each league-season. We begin by estimating  $f$  from our observational dataset of half-innings from 2010 to 2019. The response variable is a binary indicator denoting whether the pitcher’s team won the game, and the features are the inning number  $I$ , the runs allowed through that half-inning  $R$ , the league, and the season. Note that if a home team leads after the top of the 9<sup>th</sup> inning, then the bottom of the 9<sup>th</sup> is not played. Therefore, to avoid selection bias, we exclude all 9<sup>th</sup> inning instances in which a pitcher pitches at home.

With enough data, the empirical grid (e.g., binning and averaging over all combinations of  $I$  and  $R$  within a league-season) is a great estimator of  $f$ . In Figure 13a we visualize the empirical grid fit on a dataset of all half-innings from 2019 in which the home team is in the National League. The function  $f$  should be monotonic decreasing in  $R$ . In particular, as a pitcher allows more runs through a fixed number of innings, his team is less likely to win the game. It should also be monotonic increasing in  $I$  because giving up  $R$  runs through  $I$  innings is worse than giving up  $R$  runs through  $I + i$  innings for  $i > 0$ , since giving up  $R$  runs through  $I + i$  innings implies a pitcher gave up no more than  $R$  runs through  $I$  innings. The empirical grid, however, is not monotonic in either  $R$  or  $I$  because each league-season dataset is not large enough. Moreover, even when we use our entire dataset of all half-innings from 2010 to 2019, the empirical grid is still not monotonic in  $R$  or  $I$ .

To force our fitted  $f$  to be monotonic, we use XGBoost with monotonic constraints, tuned using cross validation (Chen and Guestrin, 2016). We visualize our 2019 NL XGBoost fit in Figure 13b. We indeed see that the fitted  $f$  is decreasing in  $R$  and increasing in  $I$ . Additionally,  $R \mapsto f(I, R)$  is mostly convex: if a pitcher has already allowed a high number of runs, there is a lesser relative impact of throwing an additional run on winning the game. Nonetheless, XGBoost overfits, especially towards the tails (e.g., for  $R$  large). For instance, the 2019 NL XGBoost model indicates that the probability of winning a game after allowing 10 runs through 9 innings is about 0.11, which is too large.

As there is not enough data to use machine learning to fit a separate grid for each league-season

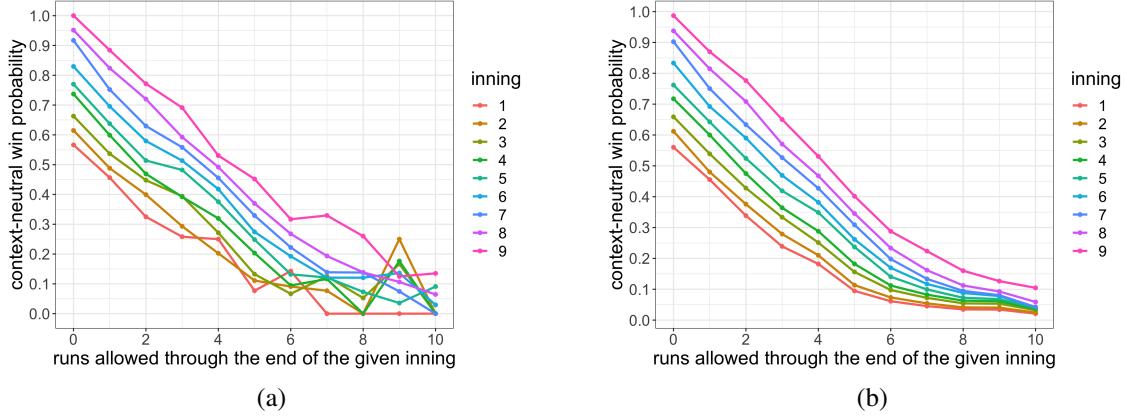


Figure 13: Estimates of the 2019 National League function  $R \mapsto f(I, R)$  using the empirical grid (left) and XGBoost with monotonic constraints (right).

without overfitting, we turn to a parametric mathematical model. Indeed, the power of parameterization is that it distills the information of a dataset into a concise form (e.g., into a few parameters), allowing us to create a strong model from limited data. Because the runs allowed in a half-inning is a natural number, we begin our parametric quest by supposing that the runs allowed in a half-inning is a  $\text{Poisson}(\lambda)$  random variable. In particular, denoting the runs allowed by the pitcher's team's batters in inning  $i$  by  $X_i$  and the runs allowed by the opposing team in inning  $i$  (for innings  $i$  after the pitcher exits the game), we assume

$$X_i, Y_i \stackrel{i.i.d.}{\sim} \text{Poisson}(\lambda). \quad (\text{A.1})$$

Then the probability that a pitcher wins the game after allowing  $R$  runs through  $I$  innings, assuming win probability in overtime is  $1/2$ , is

$$f(I, R | \lambda) := \mathbb{P}\left(\sum_{i=1}^9 X_i > R + \sum_{i=I+1}^9 Y_i\right) + \frac{1}{2} \cdot \mathbb{P}\left(\sum_{i=1}^9 X_i = R + \sum_{i=I+1}^9 Y_i\right). \quad (\text{A.2})$$

If  $I = 9$ , this is equal to

$$\mathbb{P}(\text{Poisson}(9\lambda) > R) + \frac{1}{2} \cdot \mathbb{P}(\text{Poisson}(\lambda) = R). \quad (\text{A.3})$$

If  $I < 9$ , it is equal to

$$\mathbb{P}(\text{Skellam}(9\lambda, (9-I-1)\lambda) > R) + \frac{1}{2} \cdot \mathbb{P}(\text{Skellam}(9\lambda, (9-I-1)\lambda) = R), \quad (\text{A.4})$$

noting that the Skellam distribution arises as a difference of 2 independent Poisson distributed random variables. Then, we estimate  $\lambda$  separately for each league-season by computing each team's mean runs allowed in each half inning, and then averaging over all teams.

In Figure 14a we visualize the estimated  $f$  according to our Poisson model (A.2) using the 2019 NL  $\lambda$ . We see that  $f$  is decreasing in  $R$ , increasing in  $I$ , convex in the tails of  $R$ , and is smooth. Nonetheless, some of the win probability values from this model are unrealistic. For instance, it implies the probability of winning the game after shutting out the opposing team through 9 innings is about 99%, which is too high, and the probability of winning the game after allowing 10 runs through 9 innings is about 1%, which is too low.

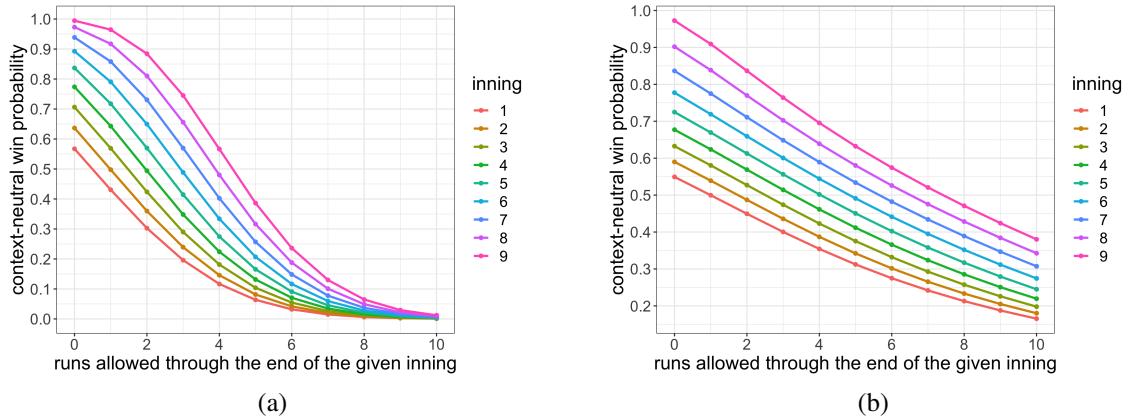


Figure 14: Estimates of the 2019 National League function  $R \mapsto f(I, R)$  using our Poisson model (A.2) with constant  $\lambda$  (left) and our Poisson model (A.8) with a truncated normal prior (A.7) on 2 team strength parameters  $\lambda_X$  and  $\lambda_Y$  (right).

The win probability values at both tails of  $R$  are too extreme in our original Poisson model (A.6) because we assume both teams have the same mean runs per inning  $\lambda$ . This is an unrealistic assumption: in real life, a baseball season involves teams of varying strength playing against each other. When teams of differing batting strength play each other, win probabilities differ. For instance, when a great hitting team allows 7 runs to a terrible hitting team, the great hitting team has a larger probability of coming back to win the game than a worse hitting team would. Thus, accounting for random differences in team strength across games should flatten the  $f(I, R)$  grid.

On this view, it is more realistic to assume the pitcher's team and the opposing team have their own runs scored per inning parameters,

$$X_i \stackrel{i.i.d.}{\sim} \text{Poisson}(\lambda_X) \quad \text{and} \quad Y_i \stackrel{i.i.d.}{\sim} \text{Poisson}(\lambda_Y), \quad (\text{A.5})$$

and

$$f(I, R | \lambda_X, \lambda_Y) := \mathbb{P} \left( \sum_{i=1}^9 X_i > R + \sum_{i=I+1}^9 Y_i \right) + \frac{1}{2} \cdot \mathbb{P} \left( \sum_{i=1}^9 X_i = R + \sum_{i=I+1}^9 Y_i \right). \quad (\text{A.6})$$

Moreover, to capture the variability in team strength across each of the 30 MLB teams, we impose a positive normal prior,

$$\lambda_X, \lambda_Y \sim \mathcal{N}_+(\lambda, \sigma_\lambda^2). \quad (\text{A.7})$$

We estimate the prior hyperparameters  $\lambda$  and  $\sigma_\lambda$  separately for each league-season by computing each team's mean and s.d. of the runs allowed in each half inning, respectively, and then averaging over all teams.

Given  $\lambda_X$  and  $\lambda_Y$ , we compute Formula (A.6) similarly as before using the Poisson and Skellam distributions. We use Monte Carlo integration with  $B = 100$  samples to estimate the posterior mean grid,

$$f(I, R | \lambda, \sigma_\lambda^2) \approx \frac{1}{B} \sum_{b=1}^B f(I, R | \lambda_X^{(b)}, \lambda_Y^{(b)}), \quad (\text{A.8})$$

where  $\lambda_X^{(b)}$  and  $\lambda_Y^{(b)}$  are i.i.d. samples from the prior distribution (A.7).

In Figure 14b we visualize the estimated  $f$  according to this Poisson model (A.8), with prior (A.7), using the 2019 NL  $\lambda$  and  $\sigma_\lambda^2$ . We see that  $f$  is mostly linear in  $R$ , rather than convex, and the values of  $f$  when  $R$  is large are highly unrealistic. For instance, this model indicates that the probability of winning the game after allowing 10 runs through 9 innings is about 38%, which is way too high. This is because our model is overdispersed, i.e. the estimated prior variance  $\sigma_\lambda^2$  is too large. For example, too large of a  $\sigma_\lambda^2$  allows  $\lambda_X$  and  $\lambda_Y$  to be very far apart, so if a pitcher allows 10 runs through 9 innings and  $\lambda_X$  is much larger than  $\lambda_Y$ , then his team will have a significant chance of coming back to win.

To resolve the overdispersion issue, we introduce a tuning parameter  $k$  designed to tune the dispersion across team strengths to match observed data,

$$\lambda_X, \lambda_Y \sim \mathcal{N}_+(\lambda, k \cdot \sigma_\lambda^2). \quad (\text{A.9})$$

In particular, we use  $k = 0.28$ , which minimizes the log-loss between the observed win/loss column and predictions from the induced grid  $f(I, R | \lambda, \sigma_\lambda^2, k)$ . In Figure 15 we visualize the estimated  $f$

according to our Poisson model (A.8), with tuned dispersion prior (A.9), using the 2019 NL  $\lambda$  and  $\sigma_\lambda^2$ . We see that  $f$  is decreasing in  $R$ , increasing in  $I$ , and convex when  $R$  is large. In particular, it looks like a smoothed version of the XGBoost grid from Figure 13b. Additionally, the values of the grid at both tails of  $R$  seem reasonable. For instance, the model indicates that allowing 0 runs through 9 innings has about a 97% win probability, which is more reasonable than before. For all of these reasons, we use this model for the grid  $f$  to compute Grid WAR for starting pitchers.

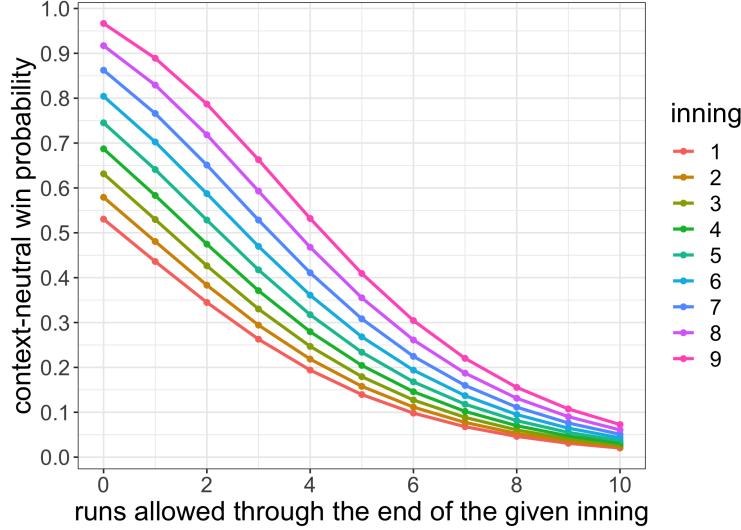


Figure 15: Estimates of the 2019 National League function  $R \mapsto f(I, R)$  using our Poisson model (A.8) with tuned dispersion prior (A.9).

## B Estimating pitcher quality using Empirical Bayes

In this Section, we describe our parametric Empirical Bayes estimators  $\hat{\mu}_p^{\text{GWAR}}$  and  $\hat{\mu}_p^{\text{FWAR}}$  of pitcher  $p$ 's quality.

### B.1 Empirical Bayes estimators of pitcher quality using Grid WAR

We begin with  $\hat{\mu}_p^{\text{GWAR}}$  which estimates pitcher quality using pitcher  $p$ 's previous games' Grid WAR and number of games played. Specifically, index each starting pitcher by  $p \in \{1, \dots, \mathcal{P}\}$  and index pitcher  $p$ 's games by  $g \in \{1, \dots, N_p\}$ . Let  $X_{pg}$  denote pitcher  $p$ 's observed Grid WAR in game

\* We have found a solid model  
of  $f = f(I, R)$  !!

Don't let the perfect be the enemy of the good

↳  $f$  is not perfect since Runs are NOT Poisson

## Remainder of the Research Process

\* Grid WAR when starting pitcher exits  
in the middle of an inning

\* W<sub>rep</sub>

\* Adjustments:

park adjustment

opposing team's batting adjustment

fielding adjustment

will cover in a  
later lecture

\* Results:

compare Grid WAR to previous WAR