# Stat 451 Project

---

Matthew Nelson, Matt Baxley, Raymond Zhao, Matthew Attanasio, Rahul Vyas

# What is a Cy Young?

Award is named after Cy Young

- Won 511 games in his career
- Best season, 1901, 33 wins, 158 Ks, 1.62 ERA
- 3 no-hitters
- Hall of Fame in 1937

Cy Young Award

"...the Cy Young Award is voted upon by the Baseball Writers' Association of America prior to the beginning of the postseason. From 1956 to 1966, the award was given to one pitcher, but has been given to one pitcher per league since 1967. The vote totals are based on a weighted points system." (MLB)

# The Data

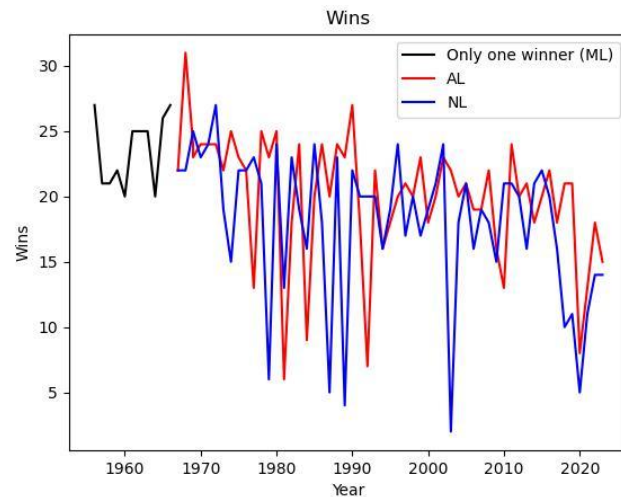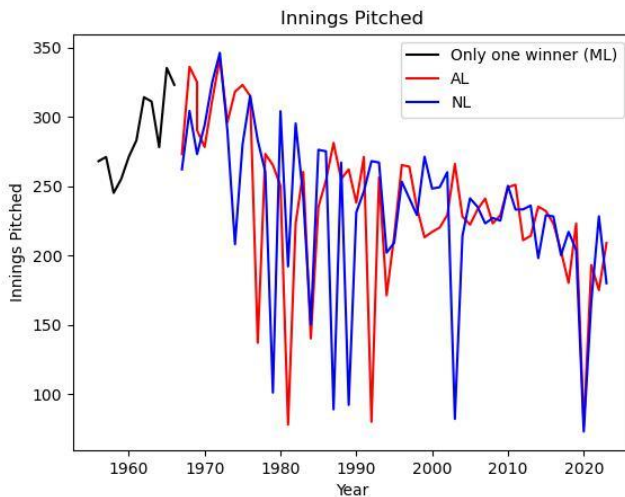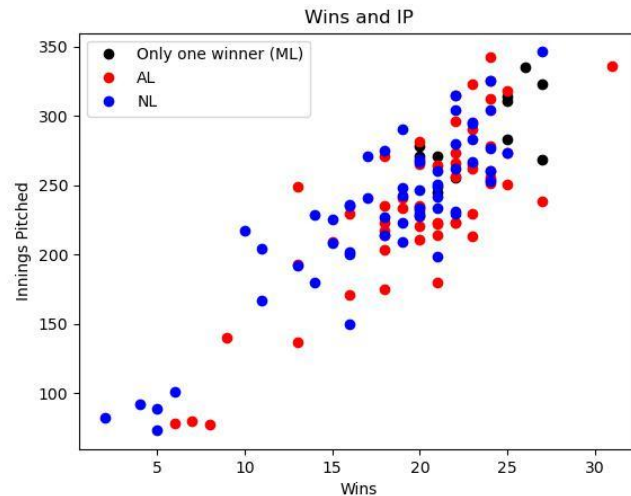2 main datasets were combined for analysis:

- All Cy Young winners since 1956
- Statistics of every pitcher in both the National League and American League (which make up the MLB) for the last 50 years (1974-2023).
  - Different types of statistics, including many advanced ones - 4 separate CSVs for each league's season (400 CSVs total)
    - Basic Statistics
    - Value
    - Advanced Statistics
    - Win Probability
  - Team Win % that year also added to dataset



All data taken from baseball-reference.com

# Correlation in the Data


Wins and IP

1. Find correlating factors
   a. W and IP: 0.85
   b. SO and WAR: 0.71
2. Plot factors over time
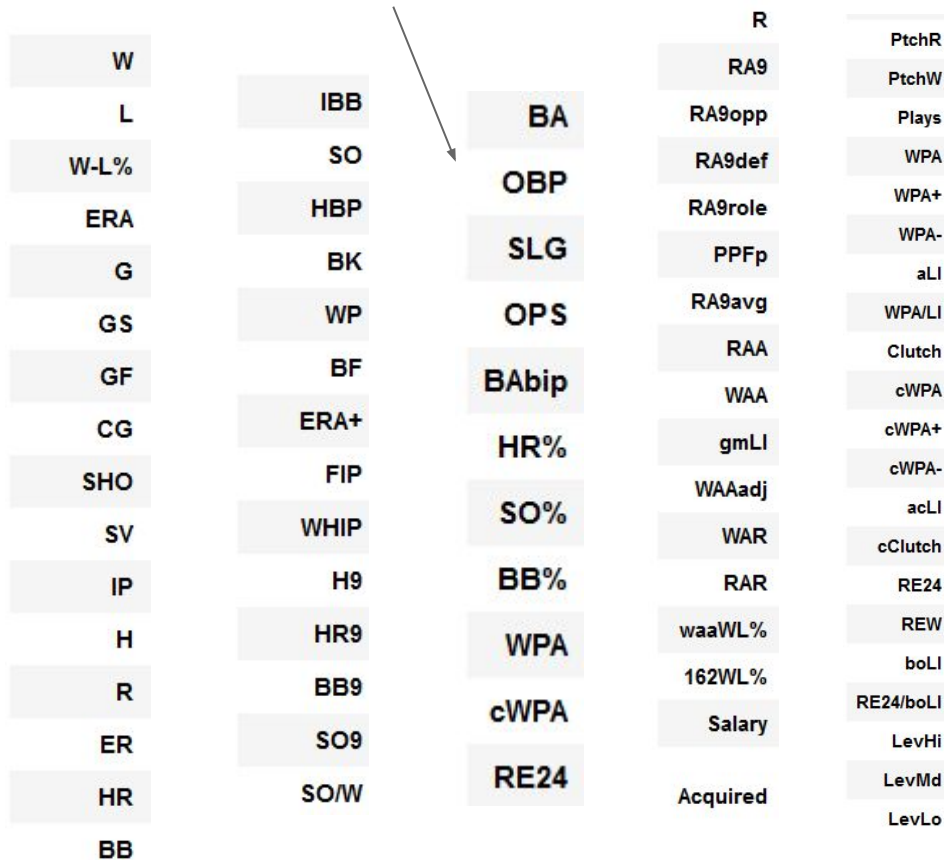

Innings Pitched


Wins

# The Plan

- Perform logistic regression on the dataset, to determine which person each year has the **highest likelihood** of being a Cy Young winner
  - Since there are many, many pitchers, we filtered the pitchers to those with 70 or more innings pitched that season.
- Since there are 2 Cy Young winners per year (1 in NL and 1 in AL), we must **oversample** the data, as it is very imbalanced
- 50/50 Train-Test split, oversample training data and test on non-oversampled testing data
- Observe model performance on entire history of MLB pitching

# What Features to Use?

- Advanced statistics on top of many basic statistics leads to a horrifying amount of features
- How do we figure out which ones to use?

(for anyone who knows baseball: these are averages of batters against the pitcher)

| | | | | |
|---|---|---|---|---|
| W | | | R | PtchR |
| L | IBB | BA | RA9 | PtchW |
| W-L% | SO | OBP | RA9opp | Plays |
| ERA | HBP | | RA9def | WPA |
| G | BK | SLG | RA9role | WPA+ |
| GS | WP | OPS | PPFp | WPA- |
| GF | BF | | RA9avg | aLI |
| CG | ERA+ | BAbip | RAA | WPA/LI |
| SHO | FIP | HR% | WAA | Clutch |
| SV | WHIP | | gmLI | cWPA |
| IP | H9 | SO% | WAAadj | cWPA+ |
| H | HR9 | BB% | WAR | cWPA- |
| R | BB9 | WPA | RAR | acLI |
| ER | SO9 | cWPA | waaWL% | cClutch |
| HR | SO/W | RE24 | 162WL% | RE24 |
| BB | | | Salary | REW |
| | | | Acquired | boLI |
| | | | | RE24/boLI |
| | | | | LevHi |
| | | | | LevMd |
| | | | | LevLo |

# Feature Selection

We selected features using precision/recall and ROC area under curve for a model fitted using only that feature, along with some knowledge about baseball.
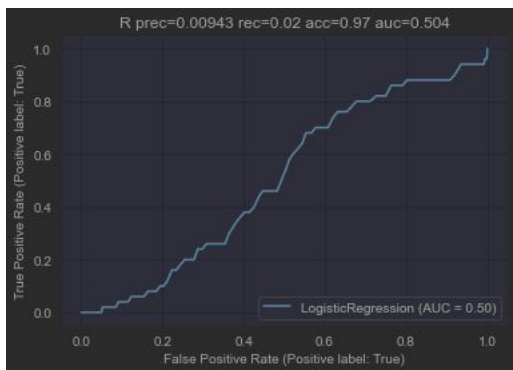
R (Runs Allowed)
Probably not getting used
Precision=0.00943, recall=0.02

WAR (Wins Over Replacement)
Probably getting used
Precision=0.291, recall=0.6

# Features Used

- **ERA** - Earned Runs Average (Average runs allowed per 9 innings pitched)
- **GS** - Games Started
- **SO** - Strikeouts
- **SV** - Saves (When the team is up by less than 3 runs and the pitcher finishes the game, that's counted as a save)
- **SHO** - Shutouts (Pitched 9 innings, opposing team scored 0 runs)
- **W** and **L** (Wins and Losses)
  - Only one win and loss is credited to a pitcher on each side each game (usually starter)
- Team Win Rate
- WAR (Wins Over Replacement)
- WPA (Win Percentage Added)
- PtchR (An advanced statistic that estimates a pitcher's contribution to their team's total runs)

# WAR

- "Wins Above Replacement"
  - How many wins a given player has added to the team above what an available replacement player would have added
- Considered one of the most important advanced statistics in baseball
- A very good measure of how much value a player is bringing to a team
- Meaning (according to baseball-reference):
  - 8+: MVP
  - 5+: All-Star
  - 2+: Starter
  - 0-2: Reserve
  - <0: Replace this player
- Similar to WPA (Win Percentage Added)
  - How much a pitcher added to their team's win chance

# What about other people's attempts?

- ESPN has a Cy Young predictor
- Citing a formula from The Neyer/James Guide to Pitchers, they use the following formula to calculate a "Cy Young" score **(highest score wins)**
  - **5/9 * (Innings Pitched) - Earned Runs + (Strikeouts / 12) + (Saves * 2.5) + Shutouts + (Wins * 6 - Losses * 2) + Victory Bonus (12 if won division, 0 otherwise)**
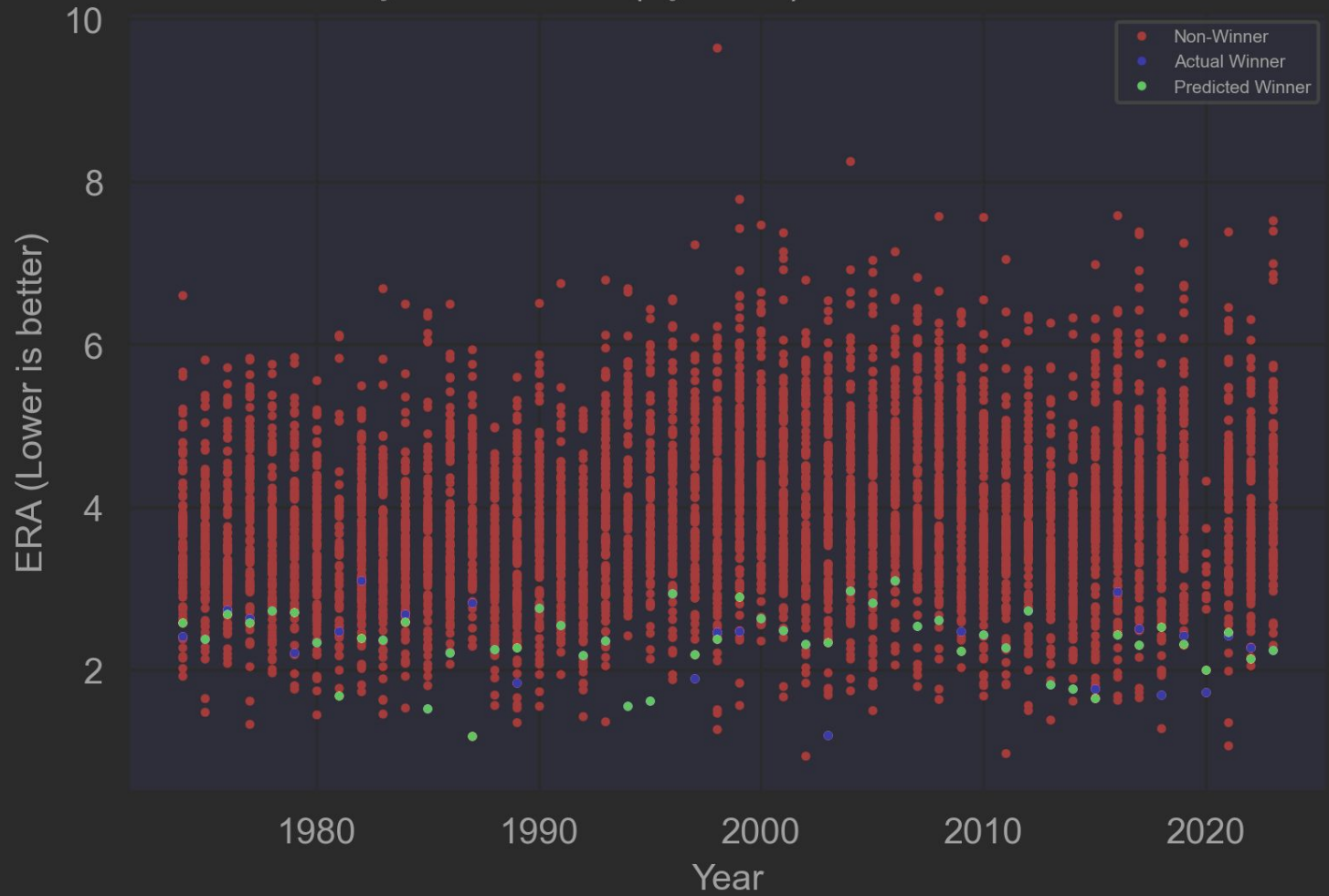
    Notably, this formula does not use any advanced statistics
- We replicate this formula, but replace the "Victory Bonus" with a bonus of 12 if they are in the top 3 teams in their league (there are 3 divisions per league)
  - No good data source available on who won each division
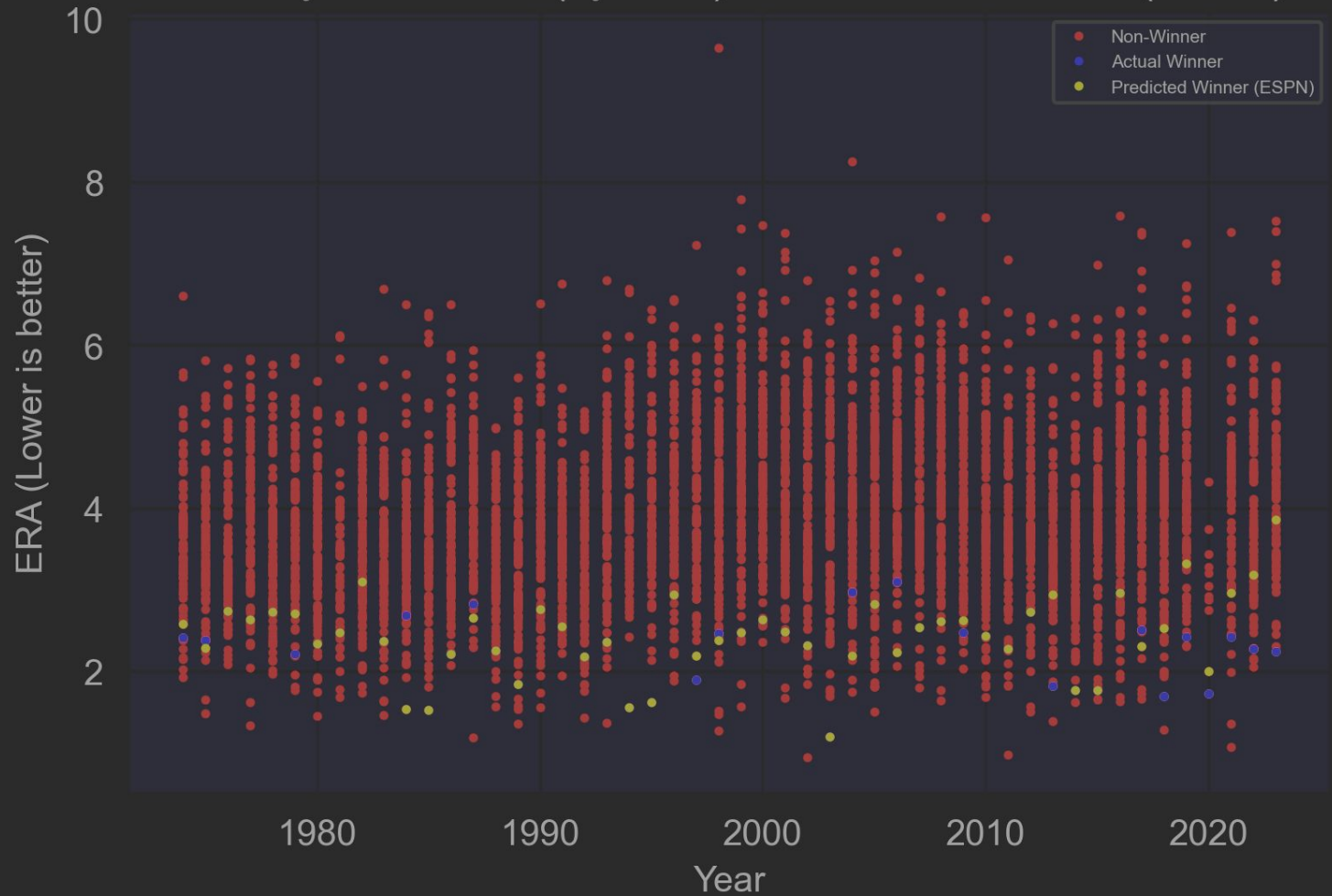- We now have 2 models to compare the performance of

NL Player Seasons (by ERA)
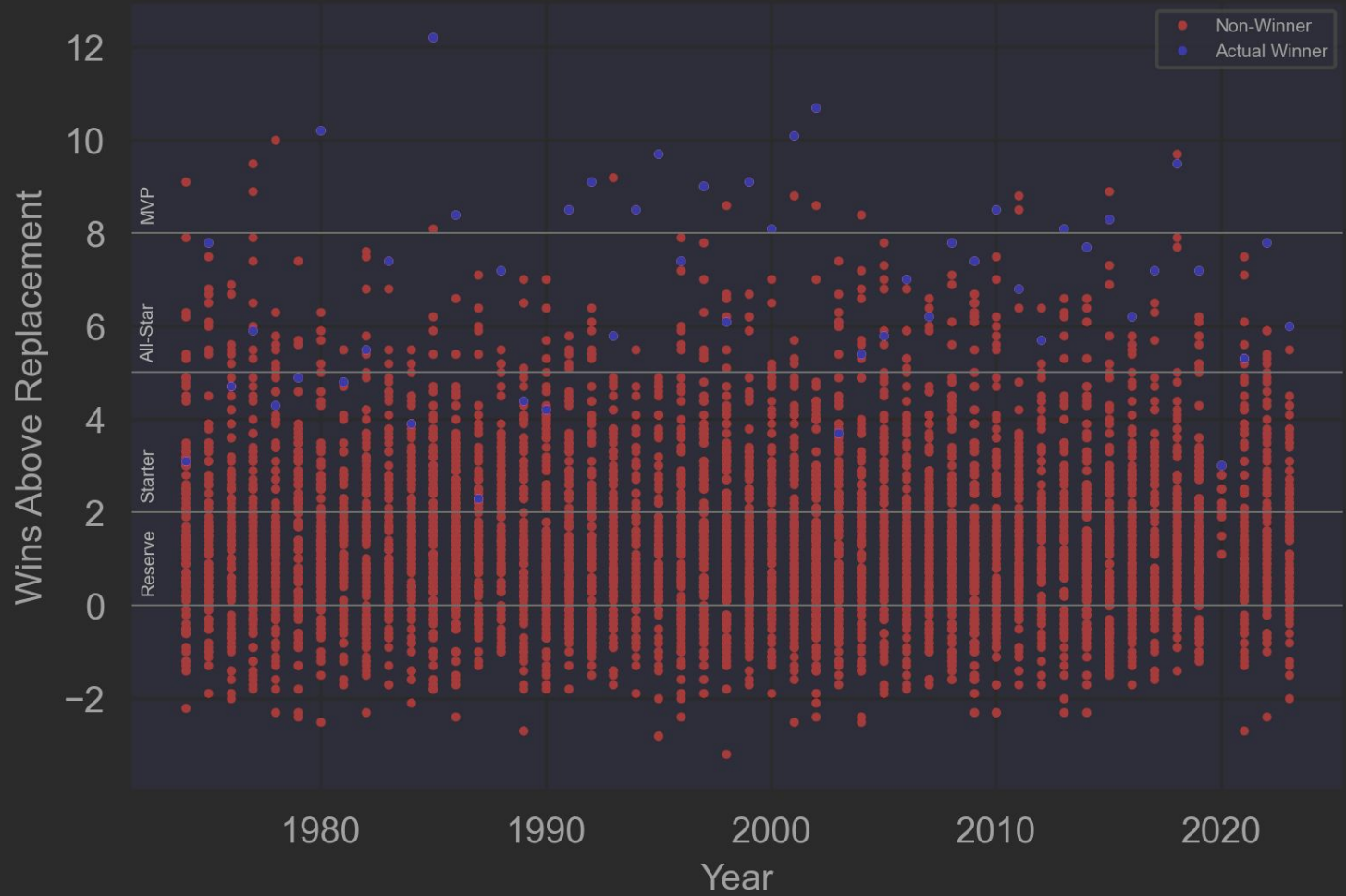
NL Player Seasons (by ERA) + Predicted Winners

NL Player Seasons (by ERA) + Predicted Winners (ESPN)

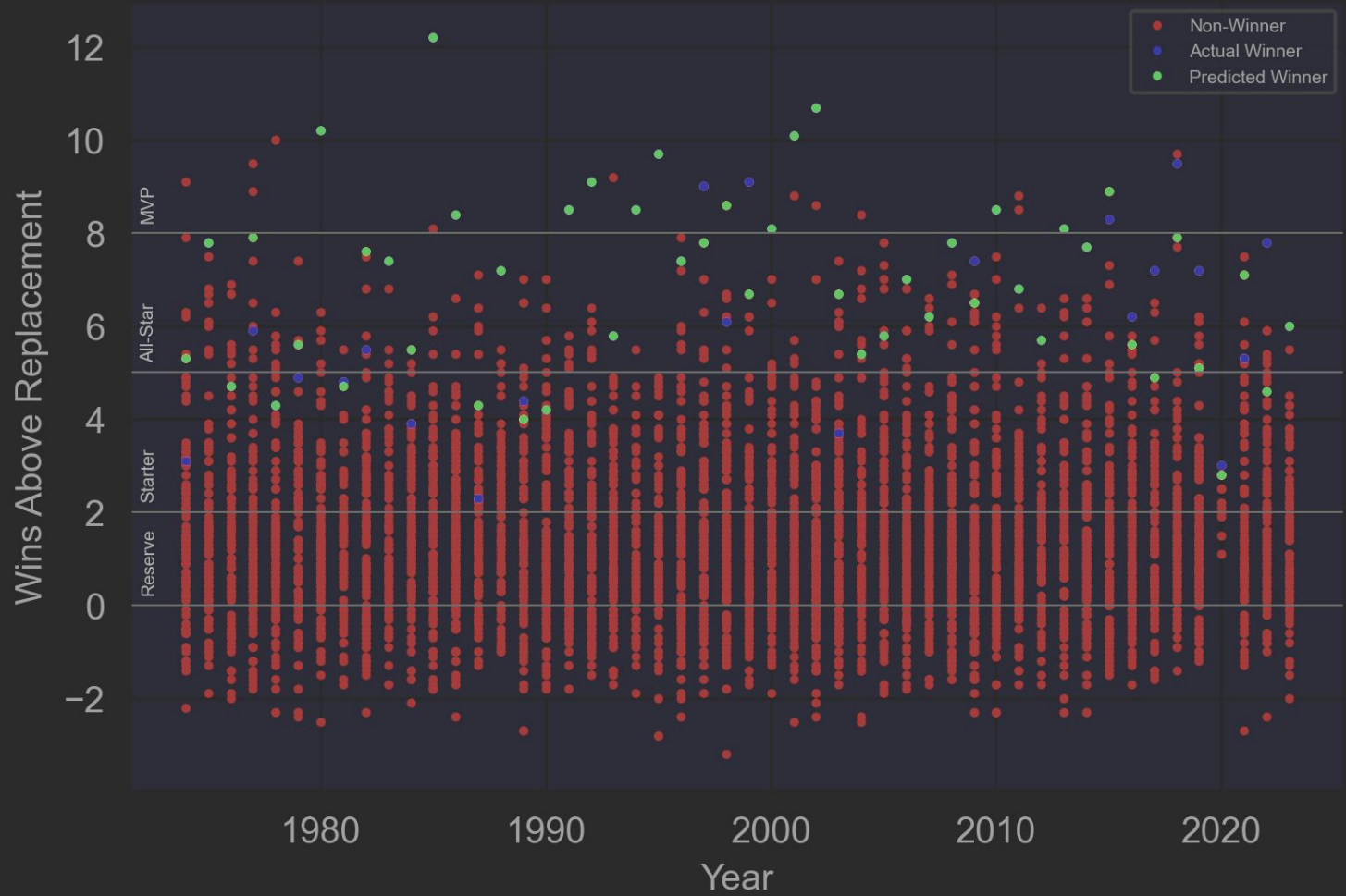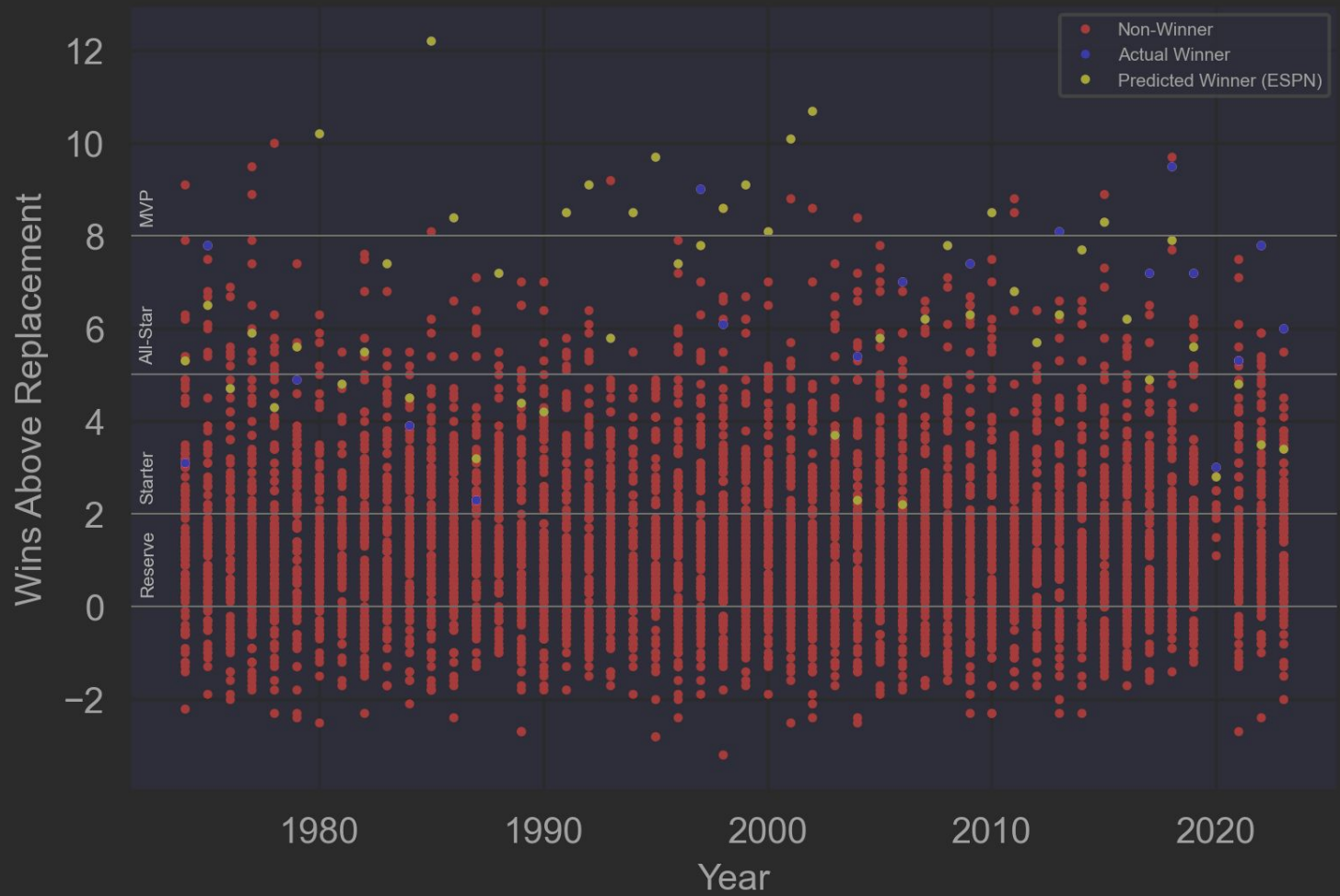NL Player Seasons (by WAR)
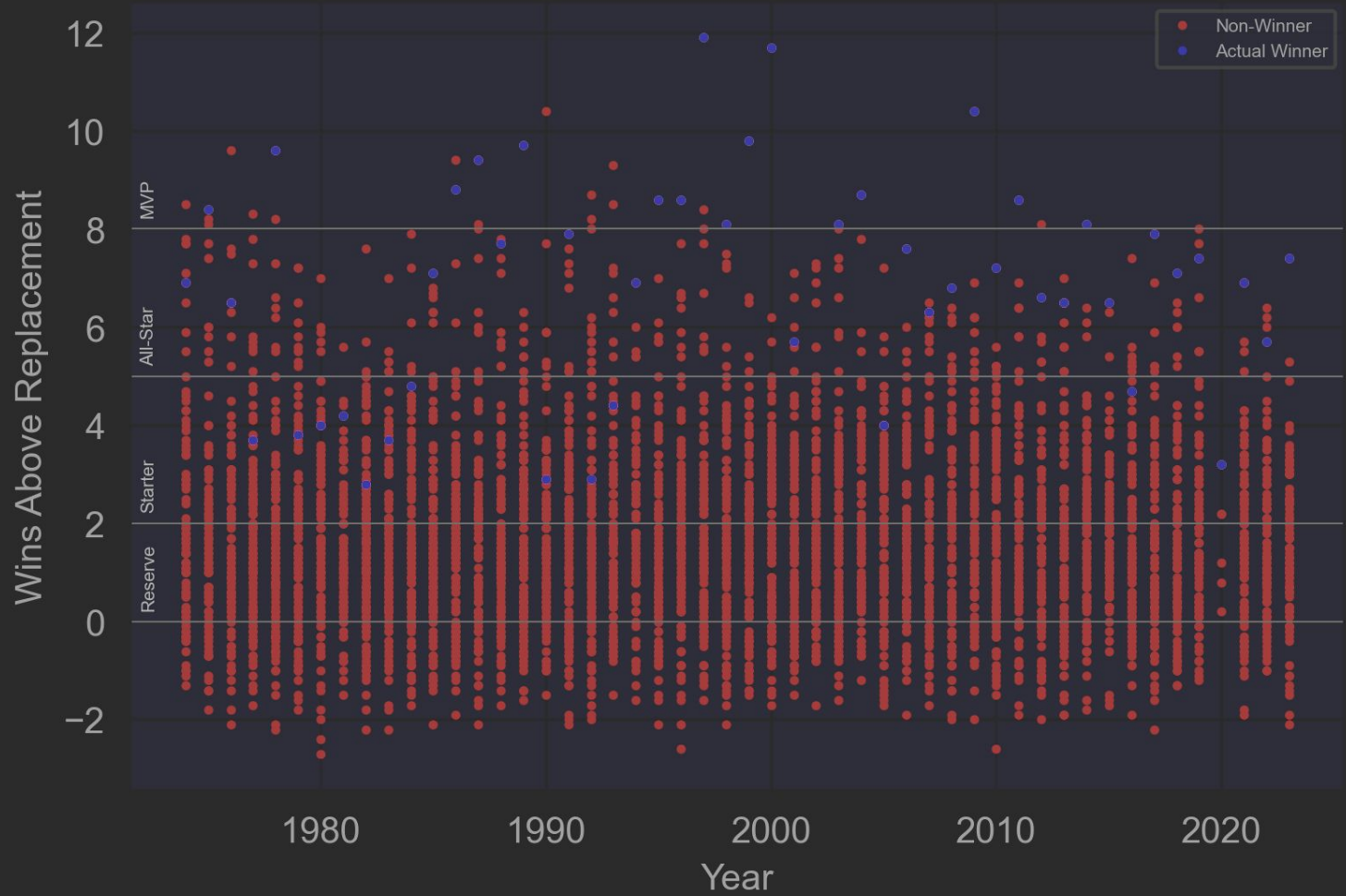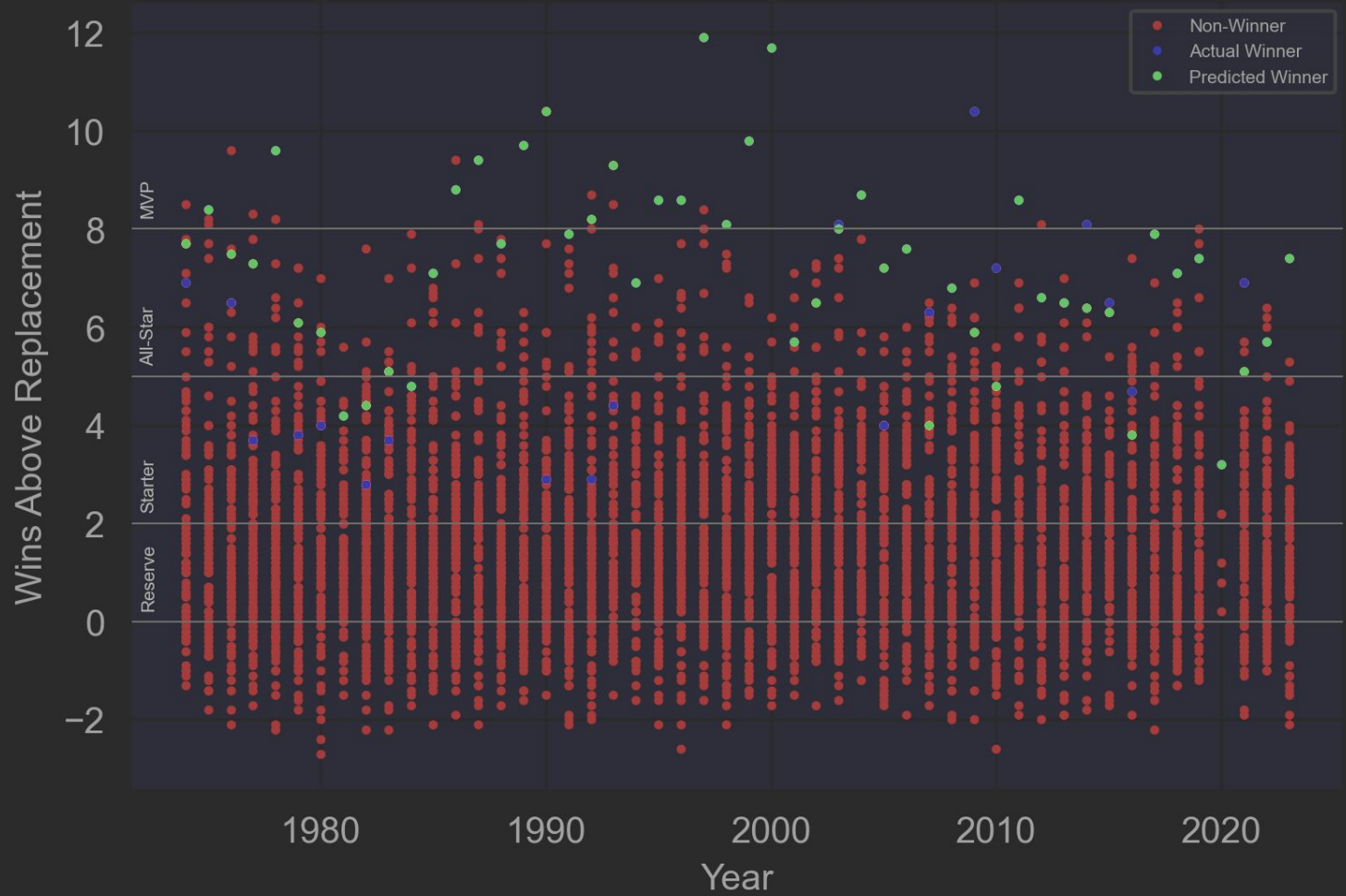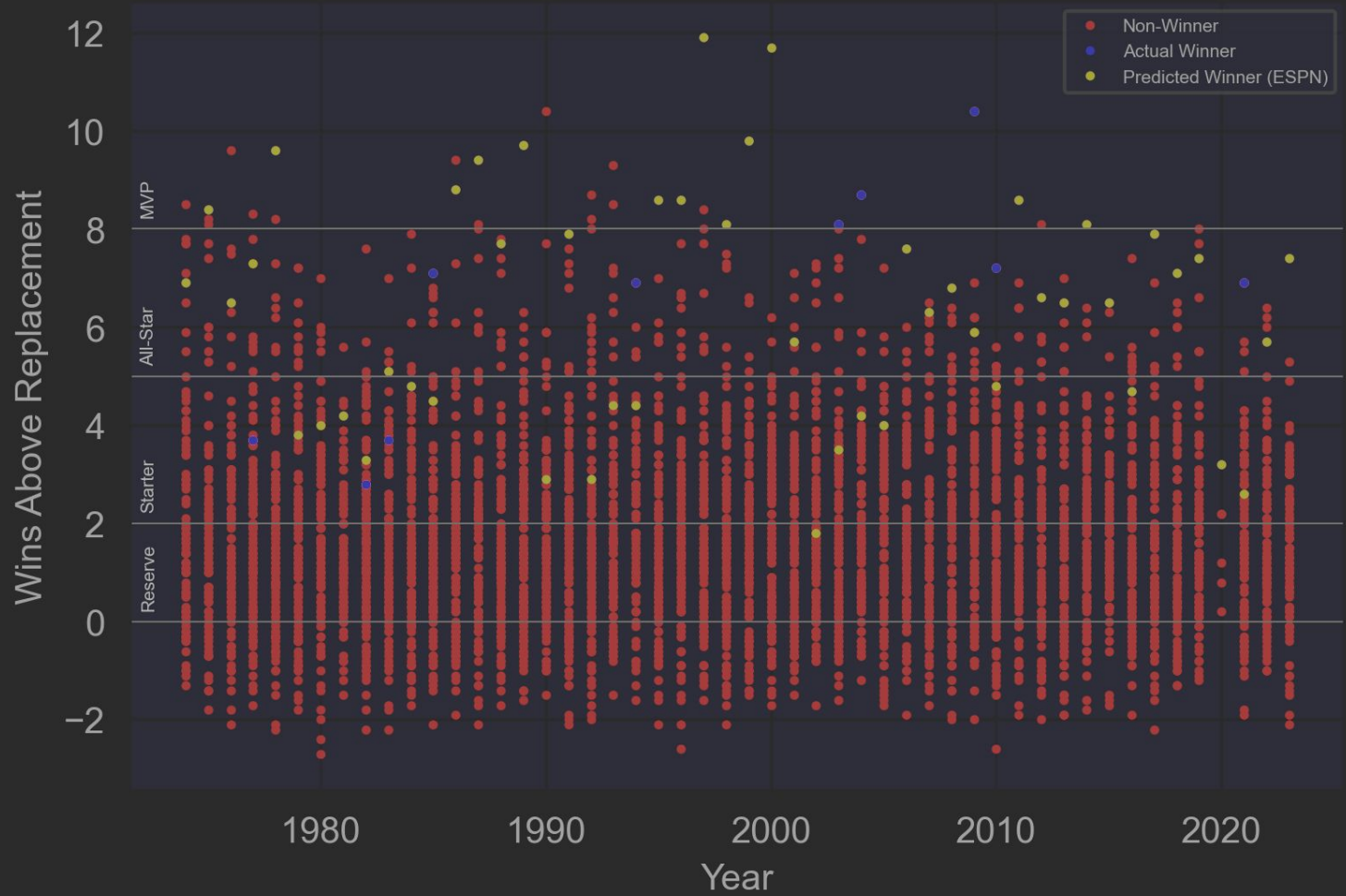
NL Player Seasons (by WAR) + Predicted Winners (ESPN)

AL Player Seasons (by WAR)

AL Player Seasons (by WAR) + Predicted Winners

AL Player Seasons (by WAR) + Predicted Winners (ESPN)

# Model Performance

Our model (on unseen test data):

- Precision: 0.36
- Recall: 0.72
- Accuracy: 0.985 (mostly meaningless because data is imbalanced)
- Area under ROC curve: 0.984
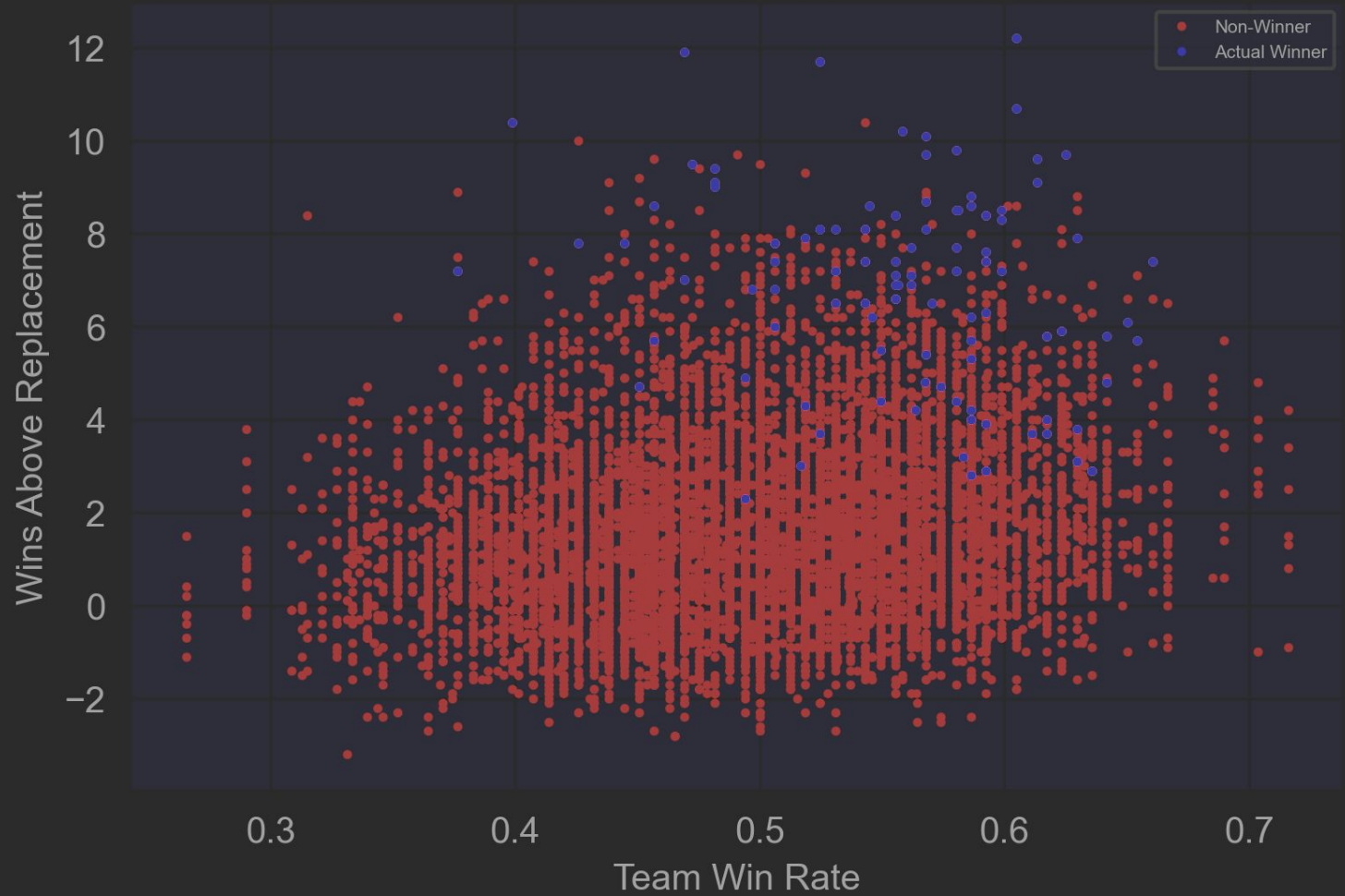
ESPN (on unseen test data):

- Precision: 0.36
- Recall: 0.72
- Accuracy: 0.985
- Area under ROC curve: 0.973
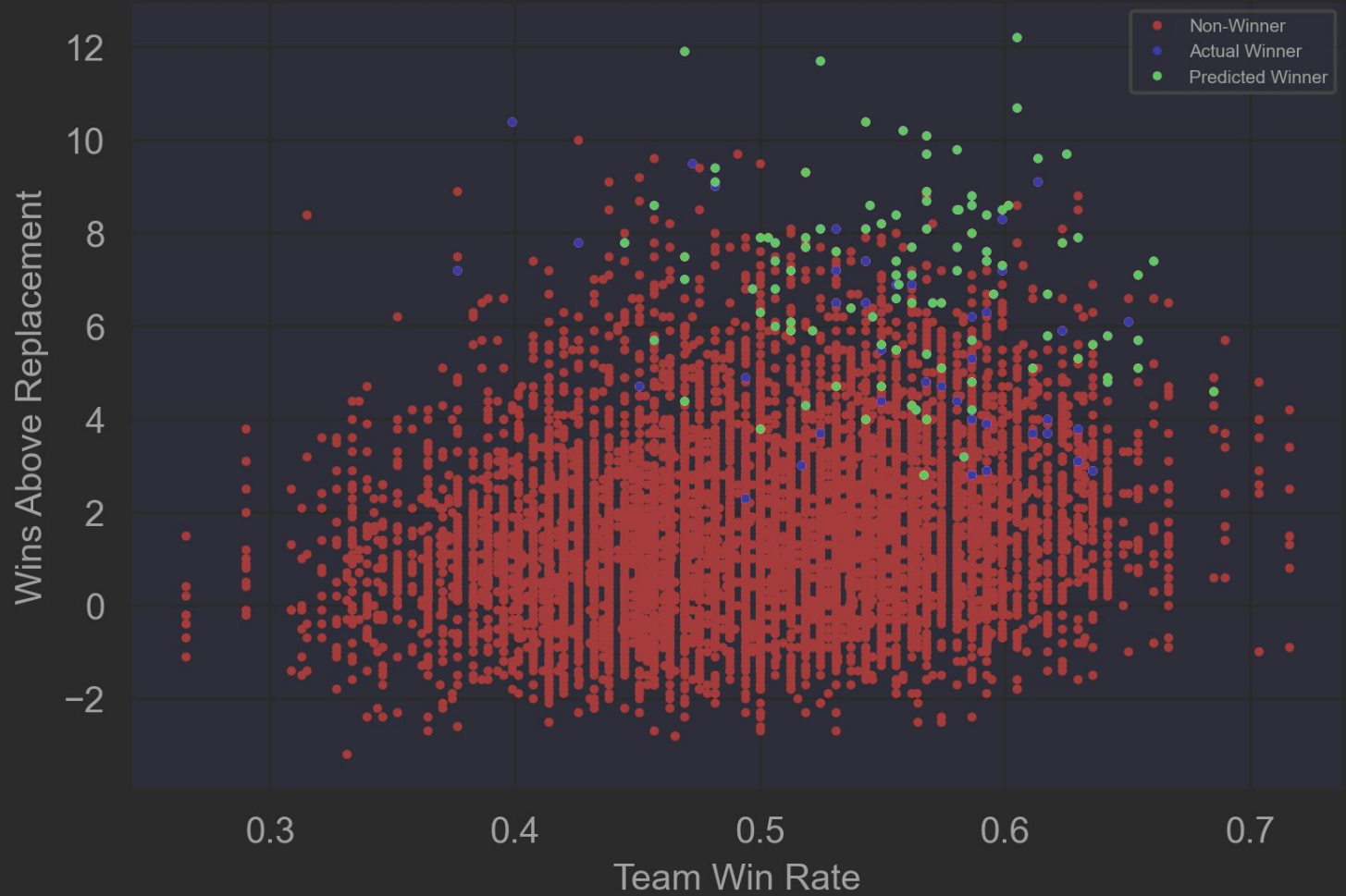
Performance was incredibly similar

**One more question:**

Does performance of the
team a pitcher is on
unfairly influence who wins
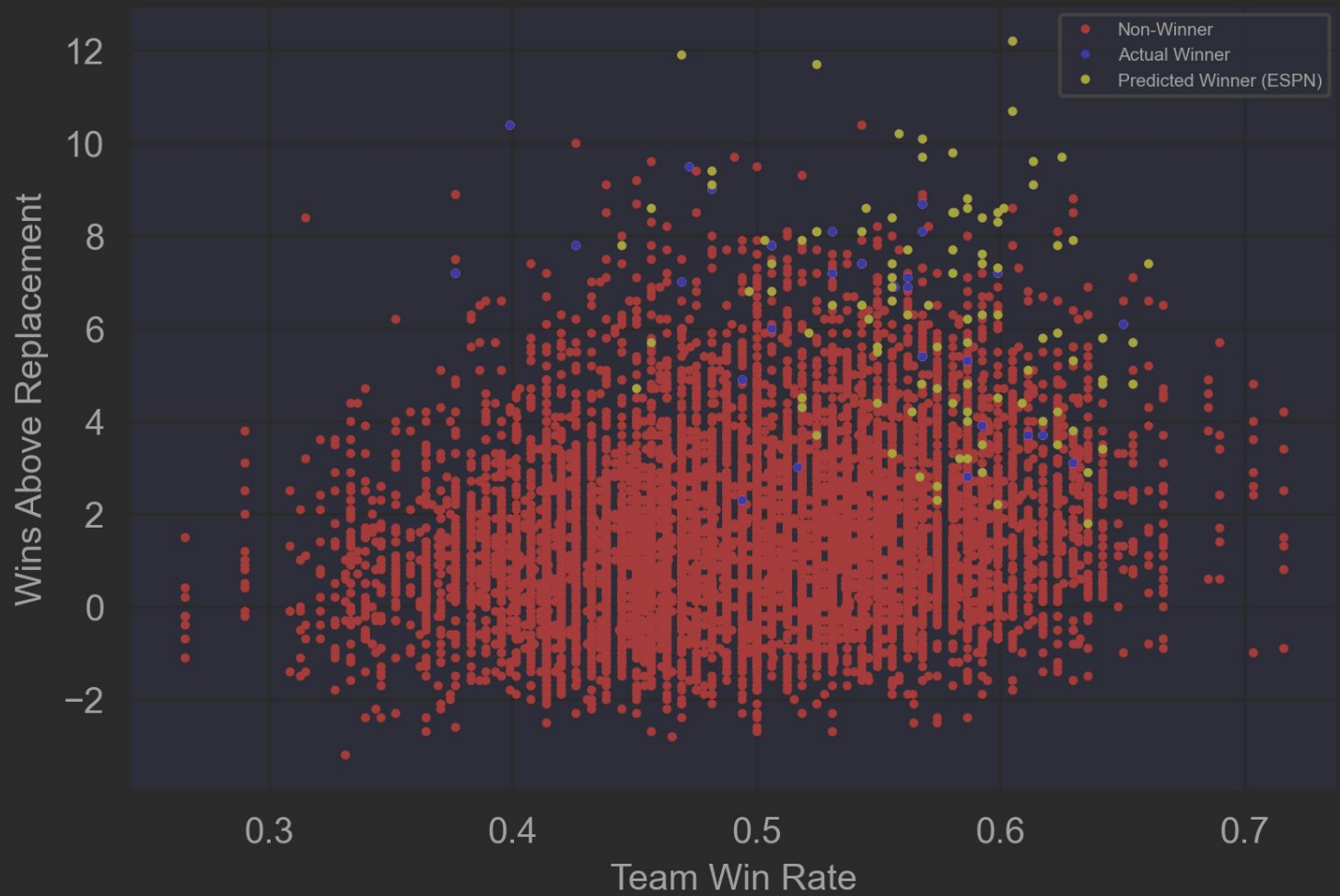the Cy Young award?

WAR vs. Team Win Rate (MLB)

WAR vs. Team Win Rate (MLB) + Predicted Winners

WAR vs. Team Win Rate (MLB) + Predicted Winners (ESPN)

# Answers:

- Cy Young Voters aren't very good at picking high-value pitchers that are on bad teams
- They like voting for pitchers that are on teams that are winning
- Since models can only reflect the results of these votes, they also both emphasize team win rate

# "Taking it Home"

- Baseball is random
  - "No two plays are the same"
- Baseball is volatile
  - Green highlight is MVP year
  - One of the most telling statistics of a season, WAR, is unpredictable between seasons
- WAR and WPA are great predictors of player impact
  - Their effectiveness at predicting who wins Cy Young is good, but limited, because voters like pitchers that are on winning teams
- Cy Young voters cannot be completely modeled using just statistics
  - There are likely intangibles that lead them to voting one way or another, that cannot be quantified through a model like this

| Year | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 |
|------|------|------|------|------|------|------|------|
| Bellinger | 3.9 | 4.2 | 8.6 | 1.5 | -1.7 | 1.2 | 4.4 |
| Yelich | 3.7 | 7.3 | 7 | 0.5 | 1.3 | 2.7 | 3.6 |