

Who's On First



Shanshan Ding
Insight Data Science

Fantasy Baseball

March

Get together with friends to each draft a team.

April to September

Everyday, decide who to play and who to bench.

October

Overall winner gets prize or bragging rights.



I hit well against Flanders, but Mr. Burns throws a lot of curveballs.

The Small Sample Size Challenge



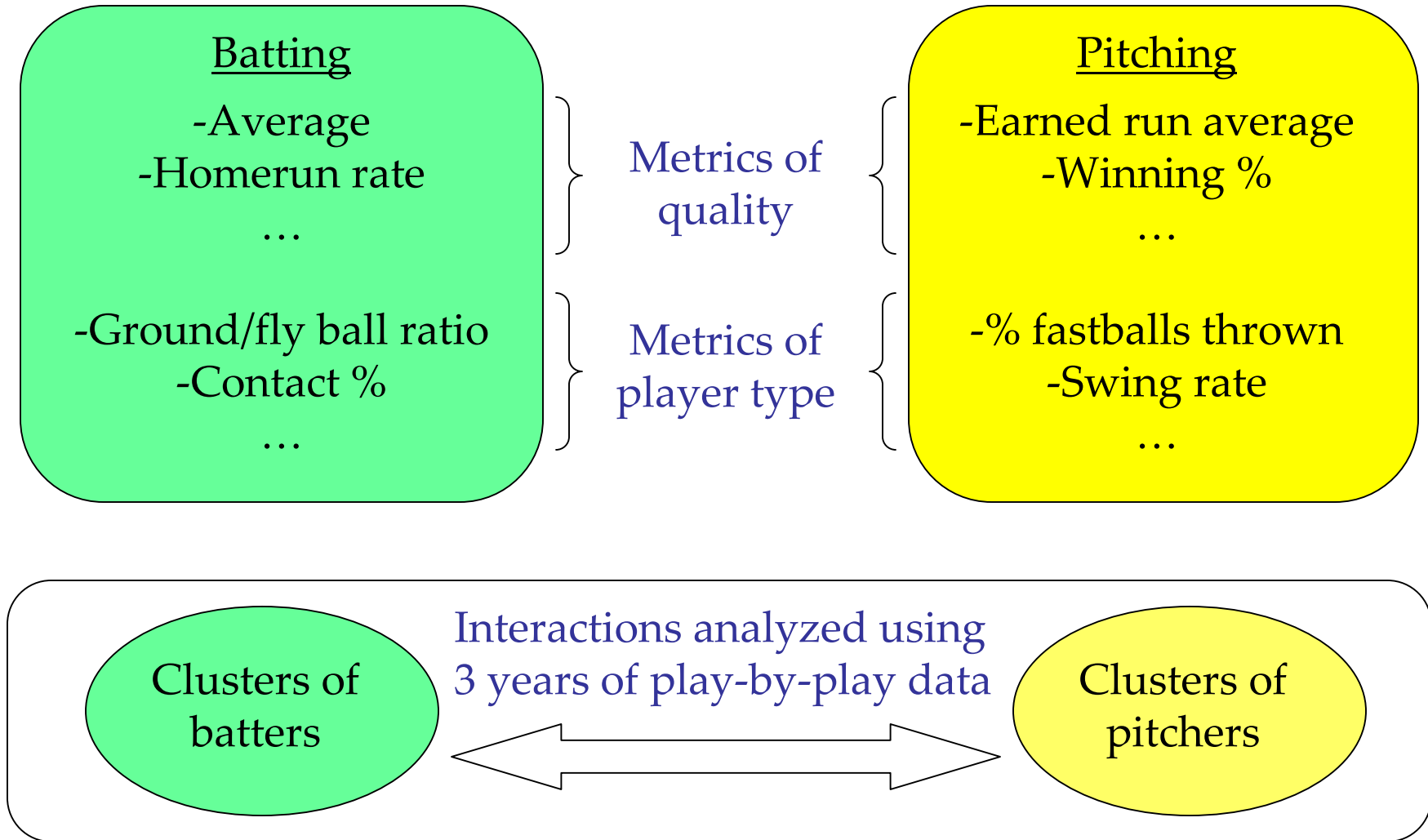
Traditional predictions: simply look up the batter-pitcher matchup history.

Hot Batter Matchups			At-Bats	AVG
<u>Quintin Berry</u> - OF	BOS vs NYY 9/14 1:05 PM	<u>CC Sabathia</u>	2-4	.500
<u>Stephen Drew</u> - SS	BOS vs NYY 9/14 1:05 PM	<u>CC Sabathia</u>	2-3	.667



- 37% of matchups involve players who have never faced each other before;
- Another 32% have histories of 5 or less at-bats.

K-Means Clustering



Validation

Retrospectively predicted the batting averages of 154 players using **individual histories**;

Predicted the batting averages of the same players using **clustered histories**;

Predictions are accurate within

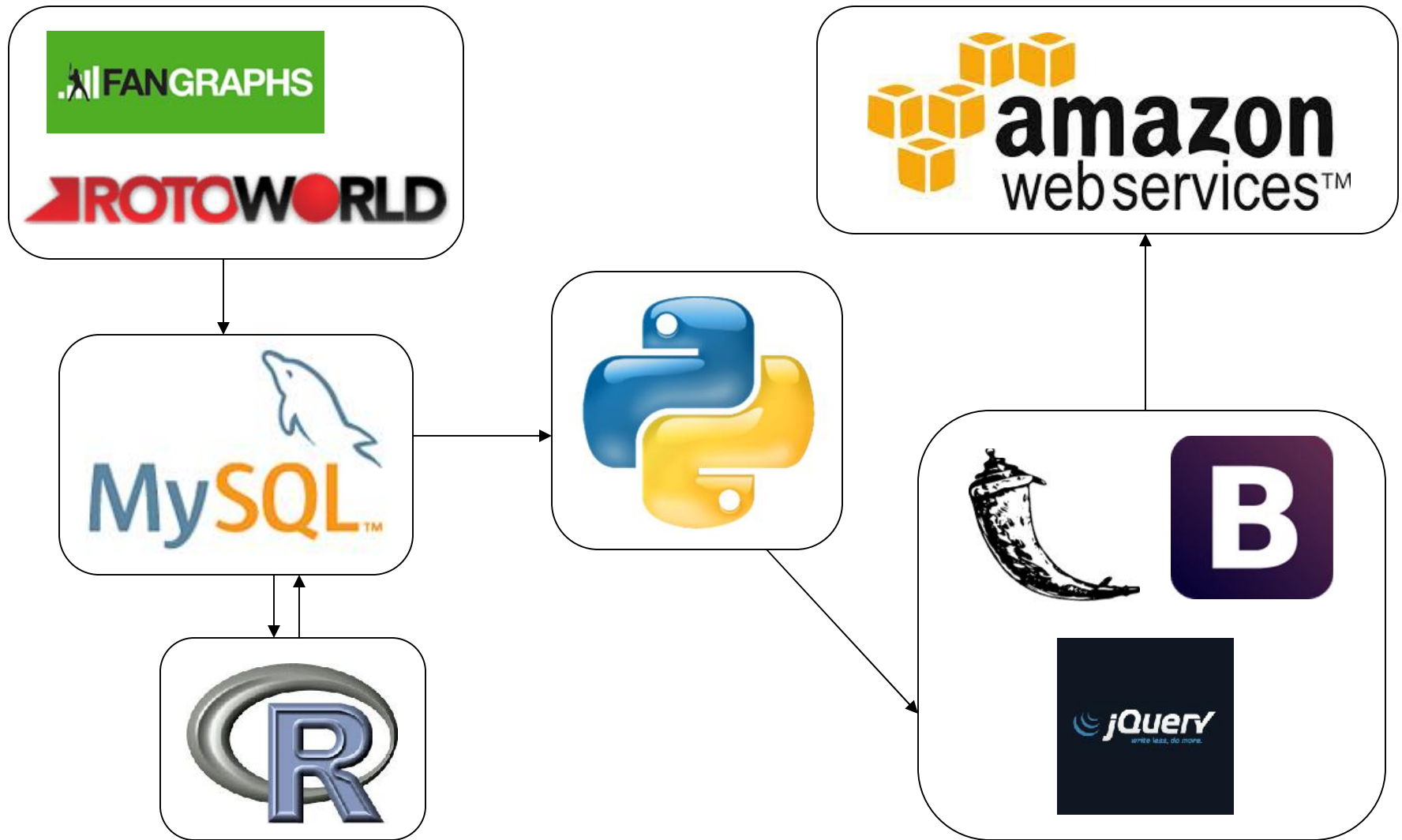
- 5% for **34** players;
- 25% for **124** players.

Predictions are accurate within

- 5% for **47** players;
- 25% for **143** players.

Compared both to the actual batting averages

Tools Used





Shanshan Ding



$$\mathbb{P}(X_t \in A | \mathcal{F}_s) = \mathbb{P}(X_t \in A | X_s)$$

