

Theoretical Derivation of the Extended James' Pythagorean Formula for Head-to-head Matchups in Baseball

Sabermetrics Modeling Group (SMALL 2024 REU)

Saad Waheed

sw21@williams.edu

Joint Work: Janine Wang, Raul Marquez, Prof. Steven J. Miller

AISC 2024 Conference
October 12th, 2024

The Original Pythagorean Win Percentage Formula

James' Pythagorean Formula

$$\text{Win Percentage} = \frac{RS^\gamma}{RS^\gamma + RA^\gamma}$$

- **RS** = Runs Scored.
- **RA** = Runs Allowed.
- **γ** : Originally taken to be 2, but empirically ≈ 1.82 for baseball.

The Original Pythagorean Win Percentage Formula

James' Pythagorean Formula

$$\text{Win Percentage} = \frac{RS^\gamma}{RS^\gamma + RA^\gamma}$$

- **RS** = Runs Scored.
- **RA** = Runs Allowed.
- **γ** : Originally taken to be 2, but empirically ≈ 1.82 for baseball.

How the Formula is Used

- **Extrapolation:** Predict a team's performance for the rest of the season based on mid-season data.
- **Evaluation:**
 - Assess if team is consistently over/under-performing.
 - Assess the value of adding players to team roster.

Applications of the Original Pythagorean Formula

Why is it Popular?

- Uses only two basic statistics [Runs Scored (RS) and Runs Allowed (RA)] but is fairly accurate despite simplicity.
- Examples:

Year	Team	RS	RA	Predicted Win % (Midseason)	Actual Win %
2009	Boston Red Sox	465	380	59.3	58.6
2011	Philadelphia Phillies	384	295	62.1	63.0
2014	San Francisco Giants	375	350	53.2	54.3
2019	New York Yankees	503	390	61.6	63.5
2023	Arizona Diamondbacks	449	422	52.9	51.9

Applications of the Original Pythagorean Formula

Why is it Popular?

- Uses only two basic statistics [Runs Scored (RS) and Runs Allowed (RA)] but is fairly accurate despite simplicity.
- Examples:

Year	Team	RS	RA	Predicted Win % (Midseason)	Actual Win %
2009	Boston Red Sox	465	380	59.3	58.6
2011	Philadelphia Phillies	384	295	62.1	63.0
2014	San Francisco Giants	375	350	53.2	54.3
2019	New York Yankees	503	390	61.6	63.5
2023	Arizona Diamondbacks	449	422	52.9	51.9

Extension to Other Sports

- Hockey ($\gamma \approx 2.1$, Dayaratna and Miller 2013)
- Basketball ($\gamma \approx 13.91$, Daryl Morey 1994)

Theoretical Basis for Pythagorean Formula (Miller 2006)

- RS and RA for a team modeled as independent random variables drawn from a continuous Weibull distributions.
- Weibull distributions chosen for their flexibility in modeling.

Theoretical Basis for Pythagorean Formula (Miller 2006)

- RS and RA for a team modeled as independent random variables drawn from a continuous Weibull distributions.
- Weibull distributions chosen for their flexibility in modeling.

Weibull Distribution

$$f_{\text{RS}}(x; \alpha, \beta, \gamma) = \begin{cases} \frac{\gamma}{\alpha} \left(\frac{x-\beta}{\alpha} \right)^{\gamma-1} \exp \left[- \left(\frac{x-\beta}{\alpha} \right)^{\gamma} \right] & \text{if } x \geq \beta \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

- α : Scale parameter (related to variance).
- β : Location parameter (translation).
- γ : Shape parameter (determines the distribution's form, whether it is more exponential or normal-like).

Shapes of the Weibull Distribution

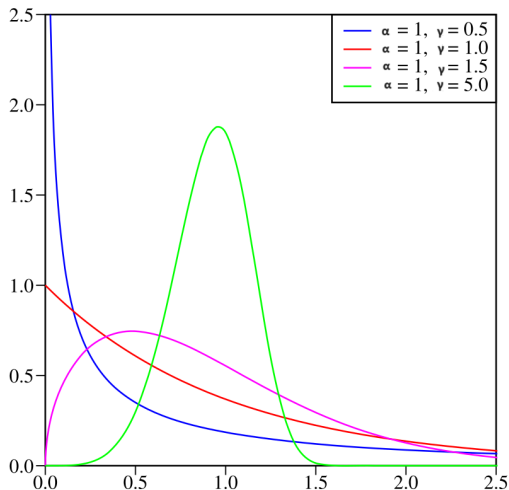


Figure: Possible Shapes of the Weibull Probability Density Function (PDF).

Weibull Distribution Expectation and Parameter Setup

Expectation of Weibull Distribution

$$\mu_{\alpha,\beta,\gamma} = \alpha\Gamma(1 + \gamma^{-1}) + \beta$$

Setting Parameters for Runs Scored (RS) and Runs Allowed (RA)

- **RS**: Modeled by a Weibull($\alpha_{RS}, \beta, \gamma$).
- **RA**: Modeled by a Weibull($\alpha_{RA}, \beta, \gamma$).
- β : Baseline for RS and RA.
- We set α_{RS} and α_{RA} so that:

$$\alpha_{RS} = \frac{RS - \beta}{\Gamma(1 + \gamma^{-1})}, \quad \alpha_{RA} = \frac{RA - \beta}{\Gamma(1 + \gamma^{-1})}.$$

Deriving the Pythagorean Win Percentage Formula

Key Steps in the Integration

$$\begin{aligned}
 \mathbb{P}(\text{Win}) &= \mathbb{P}(RS > RA) \\
 &= \int_{x=0}^{\infty} f_{RS}(x) (1 - F_{RA}(x)) dx && [\text{from Weibull PDF in (1)}] \\
 &= \int_{x=0}^{\infty} \frac{\gamma}{\alpha_{RS}} \left(\frac{x}{\alpha_{RS}} \right)^{\gamma-1} e^{-(x/\alpha_{RS})^\gamma} \left(1 - e^{-(x/\alpha_{RA})^\gamma} \right) dx \\
 &= 1 - \frac{\alpha_{RS}^\gamma}{\alpha_{RS}^\gamma + \alpha_{RA}^\gamma} \\
 &= \boxed{\frac{\alpha_{RS}^\gamma}{\alpha_{RS}^\gamma + \alpha_{RA}^\gamma}} && (\text{Pythagorean Formula})
 \end{aligned}$$

Drawbacks of the Original Pythagorean Formula

- **Poor Playoff Predictions:**
 - Less accurate in short playoff series due to small sample sizes.
- **Ignores Matchups:**
 - Assumes teams face an "average" opponent, ignoring head-to-head strengths.
- **No Home/Away Adjustment:**
 - Fails to account for home-field advantage or away-game disadvantages.

Extended Pythagorean Formula (Cleary and Miller 2020)

Adjusted Home Team Terms

$$RS_{h,adj} = RS_h \left(\frac{RA_a}{R} \right), \quad RA_{h,adj} = RA_h \left(\frac{RS_a}{R} \right)$$

RS_h (**RA_h**): Runs scored (allowed) by home team.

RS_a (**RA_a**): Runs scored (allowed) by away team.

R: League average runs scored/conceded per game.

⇒ Adjusted by a scaling factor of away team performance relative to league averages.

Extended Pythagorean Formula (Cleary and Miller 2020)

Adjusted Home Team Terms

$$RS_{h,adj} = RS_h \left(\frac{RA_a}{R} \right), \quad RA_{h,adj} = RA_h \left(\frac{RS_a}{R} \right)$$

RS_h (**RA_h**): Runs scored (allowed) by home team.

RS_a (**RA_a**): Runs scored (allowed) by away team.

R: League average runs scored/conceded per game.

⇒ Adjusted by a scaling factor of away team performance relative to league averages.

Extended Win Probability Formula

$$\begin{aligned} \mathbb{P}(\text{Home Team Wins}) &= \frac{(RS_{h,adj})^\gamma}{(RS_{h,adj})^\gamma + (RA_{h,adj})^\gamma} \\ &= \frac{(RS_h RA_a)^\gamma}{(RS_h RA_a)^\gamma + (RA_h RS_a)^\gamma} \end{aligned}$$

Performance Against Playoff Data (2001-2019)

Results

- **Extended Pythagorean Predictions:**
 - Higher seed predicted to win **80.18** series, lose **68.82**.

Performance Against Playoff Data (2001-2019)

Results

- **Extended Pythagorean Predictions:**
 - Higher seed predicted to win **80.18** series, lose **68.82**.
- **Observed Results:**
 - Higher seed won **80** series, lost **69**.

Conclusion

- Provides better predictions for head-to-head matchups than the original Pythagorean formula in playoffs, while maintaining simplicity.

Summer Research Summary

Goals for the Summer

- Derive an extended Pythagorean win-loss formula using a theoretical framework similar to Miller (2006).
- Improve formula further by incorporating team-specific factors into head-to-head matchups.

Summer Research Summary

Goals for the Summer

- Derive an extended Pythagorean win-loss formula using a theoretical framework similar to Miller (2006).
- Improve formula further by incorporating team-specific factors into head-to-head matchups.

Approach

- Treat RS_a, RA_a for the away team as the observed value of a random variable with

$$\mathbb{E}[RS_a] = \frac{n(mR) - RS_h}{n - 1} \approx mR$$

where n = # of teams in league, m = # of matches so far in season.

- Use **conditional distributions** to model Runs Scored (RS) and Runs Allowed (RA) by home team, adjusted by the away team metrics RS_a, RA_a .

Method of Conditional Distributions

Conditional Weibull Distribution for RS_h

Denote X_h as the random variables for runs scored by the home team in a fixture. Then,

$$X_h \sim \text{Weibull}(\alpha_{RS_h,adj}, \beta, \gamma).$$

- X_h dependent only on RA_a and not RS_a .
- As before, assume RS_h and RA_h to be independent.

$$f_{X_h|RA_a}(x) = \begin{cases} \frac{\gamma}{\alpha_{RS_h,adj}} \left(\frac{x - \beta}{\alpha_{RS_h,adj}} \right)^{\gamma-1} \exp \left(- \left(\frac{x - \beta}{\alpha_{RS_h,adj}} \right)^\gamma \right) & \text{if } x \geq \beta \\ 0 & \text{if } x < \beta \end{cases}$$

where

$$\alpha_{RS_h,adj} = \frac{RS_{h,adj} - \beta}{\Gamma(1 + \gamma^{-1})} = \frac{RS_h \cdot \frac{RA_a}{mR} - \beta}{\Gamma(1 + \gamma^{-1})}.$$

Expectation of X_h and Unconditional Distribution

Expectation of X_h

$$\begin{aligned}\implies \mathbb{E}[X_h \mid RA_a] &= RS_h \cdot \frac{RA_a}{mR} \\ \therefore \mathbb{E}[X_h] &= \mathbb{E}_{RA_a}[E[X_h \mid RA_a]] \\ &= \mathbb{E}_{RA_a} \left[RS_h \cdot \frac{RA_a}{mR} \right] \\ &= RS_h \cdot \frac{\mathbb{E}[RA_a]}{mR} \\ &= RS_h\end{aligned}$$

Expectation of X_h and Unconditional Distribution

Expectation of X_h

$$\begin{aligned}
 \Rightarrow \mathbb{E}[X_h \mid RA_a] &= RS_h \cdot \frac{RA_a}{mR} \\
 \therefore \mathbb{E}[X_h] &= \mathbb{E}_{RA_a}[E[X_h \mid RA_a]] \\
 &= \mathbb{E}_{RA_a} \left[RS_h \cdot \frac{RA_a}{mR} \right] \\
 &= RS_h \cdot \frac{\mathbb{E}[RA_a]}{mR} \\
 &= RS_h
 \end{aligned}$$

Unconditional Weibull Distribution for X_h

$$f_{X_h}(x) = \begin{cases} \frac{\gamma}{\alpha_{RS_h}} \left(\frac{x - \beta}{\alpha_{RS_h}} \right)^{\gamma-1} \exp \left(- \left(\frac{x - \beta}{\alpha_{RS_h}} \right)^{\gamma} \right), & \text{if } x \geq \beta \\ 0 & \text{if } x < \beta. \end{cases}$$

Deriving the Extended Pythagorean Formula

Extended Pythagorean Formula

$$\begin{aligned}
 \mathbb{P}(\text{Win}) &= \mathbb{P}(X_h > Y_h) \\
 &= \int_{\beta}^{\infty} f_{X_h|\text{RA}_a}(x) [1 - F_{Y_a}(x)] dx \\
 &= \int_{\beta}^{\infty} \frac{\gamma}{\alpha_{\text{RS}_{h,adj}}} \left(\frac{x - \beta}{\alpha_{\text{RS}_{h,adj}}} \right)^{\gamma-1} e^{-\left(\frac{x - \beta}{\alpha_{\text{RS}_{h,adj}}} \right)^{\gamma}} \left[1 - e^{-\left(\frac{x - \beta}{\alpha_{\text{RA}_a}} \right)^{\gamma}} \right] dx \\
 &= \frac{\alpha_{\text{RS}_{h,adj}}^{\gamma}}{\alpha_{\text{RS}_{h,adj}}^{\gamma} + \alpha_{\text{RA}_a}^{\gamma}} = \frac{(\text{RS}_h \cdot \frac{\text{RA}_a}{mR})^{\gamma}}{(\text{RS}_h \cdot \frac{\text{RA}_a}{mR})^{\gamma} + (\text{RA}_h \cdot \frac{\text{RS}_a}{mR})^{\gamma}} \\
 &= \boxed{\frac{(\text{RS}_h \cdot \text{RA}_a)^{\gamma}}{(\text{RS}_h \cdot \text{RA}_a)^{\gamma} + (\text{RA}_h \cdot \text{RS}_a)^{\gamma}}} \quad (\text{Extended Pythagorean Formula})
 \end{aligned}$$

Future Work

- **Analyzing Left and Right-Handed Matchups:**
 - Study the impact of batter and pitcher handedness on game outcomes.
 - Investigate how handedness affects runs scored and allowed.
 - Integrate handedness statistics into the extended Pythagorean formula.
- **Model Validation and Testing:**
 - Apply the extended Pythagorean formula to more recent data (post-2020) and compare predictive accuracy with existing metrics.

Acknowledgements

- I would like to thank Professor Steven J. Miller for his invaluable guidance and support throughout this research.
- This work was supported by Williams College and the National Science Foundation through grant DMS-2241623.

Bibliography

- ¹K. D. Dayaratna and S. J. Miller, *The pythagorean won-loss formula and hockey: a statistical justification for using the classic baseball formula as an evaluative tool in hockey*, 2013.
- ²S. J. Miller, *A Derivation of the Pythagorean Won-Loss Formula in Baseball*, 2006.
- ³R. Cleary, J. Jeffries, C. Miller, S. J. Miller, J. Murray, and N. Skiera, *Extending James' Pythagorean Formula for Head-to-head Matchups*, 2020.
- ⁴J. Dewan and D. Zminda, *STATS Basketball Scoreboard, 1993-94*, p. 17 (STATS, Inc., Oct. 1993).