

Lab: Simple Linear Regression

1. Pythagorean Win Percentage

We are given a dataset of team-seasons from 2017 to 2021,

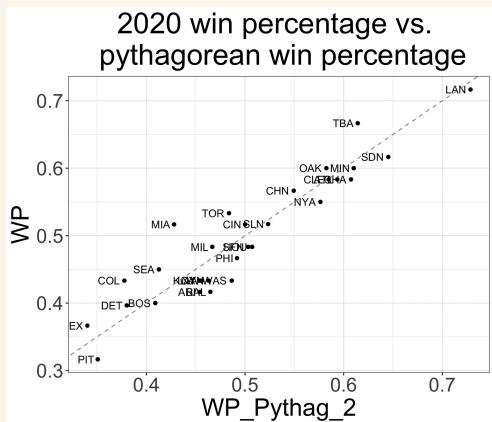
$$\left\{ \begin{array}{l} i = \text{index of } i^{\text{th}} \text{ team-season} \\ RS_i = \text{runs scored} \\ RA_i = \text{runs allowed} \\ WP_i = \text{win percentage} \end{array} \right.$$

and we want to predict end-of-season win percentage from runs scored and runs allowed. A team's deviation from this prediction is a measure of how lucky the team was.

Bill James, godfather of Sabermetrics (baseball analytics) and sports analytics, created Pythagorean Win percentage

$$\widehat{WP} = \frac{RS^2}{RS^2 + RA^2}$$

He made it up and it works quite well!



The pythagorean exponent is quite good, but is arbitrary.

- Use linear regression to find an exponent α so that the Pythagorean win percentage

$$\widehat{WP} = \frac{RS^\alpha}{RS^\alpha + RA^\alpha} \text{ best fits the data.}$$

You'll need to transform this equation to be linear in α .

Hint: divide top and bottom by RS^α

- Create a visualization to show that

$$\widehat{WP} = \frac{RS^\alpha}{RS^\alpha + RA^\alpha} \text{ is better than } \widehat{WP} = \frac{RS^2}{RS^2 + RA^2}.$$

2. Evaluating MLB general managers

We are given a dataset of MLB team payrolls and results for each season 1998–2023,

$$\left\{ \begin{array}{l} \text{row } i \leftarrow i^{\text{th}} \text{ team-season} \\ \text{win percentage} \\ \text{payroll}/\text{median payroll} \\ \text{Log}(\text{payroll}/\text{median payroll}) \end{array} \right.$$

We want to analyze the relationship between payroll and winning to evaluate general managers.

(try using Chat GPT Data Analyst for this question)

- Plot payroll/median against winning percentage.
Mark the Oakland A's and NY Yankees dots.
Remove 2020.
Add the regression line of WP as a function of
Payroll/median. Add the regression line of WP as a
function of $\log(\text{payroll}/\text{median})$.

- Now for each team-season calculate the difference between the actual WP and predicted WP using payroll/median and then $\log(\text{payroll})/\text{median}$. Add this column to the dataset. Find the average difference for each team and make a graph ordered from highest to lowest. (one graph for each model). Change the Y axis scale to wins by multiplying by 162.