

Lab: Multivariable Linear Regression

1. Four Factors in Basketball

Basketball has many dimensions:
shooting accuracy, defense, steals, blocking, rebounding...

Which correlates with (or causes) winning?

Which matter most? Least?

Dean Oliver defined four factors :

Scoring, crashing, protecting, attacking, which correspond to
shooting accuracy, rebounding efficiency, turnover rates, and free throws.

A big innovation: express these in percentage terms
(this normalizes for pace of play).

Your task is to use multivariable regression to
assess the relative impacts of the four factors on
winning. Let

$$\left\{ \begin{array}{l} y = \text{team wins in a season (outcome variable)} \\ x_1 = \text{EFG\%} - \text{Opp EFG\%} \\ x_2 = \text{TOV\%} - \text{Opp TOV\%} \\ x_3 = \text{OREB\%} - \text{DREB\%} \\ x_4 = \text{FTRate} - \text{Opp FTRate} \end{array} \right.$$

Here, $EFG\%$ = effective field goal percentage = $\frac{FG + \frac{1}{2}(3PT)}{FGA}$
 ↳ weight successful three pointers 50% more

$TOV\%$ = turnover percentage

$OREB\%$ = off. rebounds %

$DREB\%$ = def. rebounds %

$FT\%$ = free throw rate.

- First, get to know your data.
 Find means, SDs, and correlations b/t the variables
 and consider their marginal distributions.
- Fit a multivariable regression model

$$E y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4.$$
 Fit a second regression model after standardizing
 each of the x variables (to have mean 0 and sd $\frac{1}{2}$).
 Which regression tells you about the relative value
 of each of the four factors to winning?
 Order the factors by importance to winning.
 Which has better predictive performance and why
 (use math; no need to code up out-of-sample prediction calculation if you don't want to)

2. Expected outcome of a punt

We have a dataset consisting of punts,

$\left\{ \begin{array}{l} \text{row} = \text{a punt} \\ i = \text{index of } i^{\text{th}} \text{ punt} \\ y_i = \text{outcome (next yardline, from opponent's perspective) of the punt} \\ \text{ydl}_i = \text{yardline (yards from opp. endzone) of } i^{\text{th}} \text{ punt} \\ \text{pq}_i = \text{punter quality of the } i^{\text{th}} \text{ punter (I made this variable)} \\ \text{punter}_i = \text{name of punter} \end{array} \right.$

- use multivariable regression to model outcome (next yardline) as a function of yardline and punter quality.

Consider linear terms $\beta_1 \cdot \text{ydl}_i + \beta_2 \cdot \text{pq}_i$
and also transformations (e.g. quadratic, cubic, splines).

Select a model (e.g., by out-of-sample predictive performance) and visualize that model (e.g. plot expected next yardline versus yardline for various punter quality values (color)).

- Rank the punters by PYOE (punt yards over expected) and visualize the rankings. (e.g. plot punter vs. career PYOE)