

The Power of Fake Data (Priors)

Q Suppose the Dodgers have won W and lost L games thus far in the season.

How would you predict their end of season win percentage WP?

- 162 total games in the season
- no access to their schedule (e.g., ignore strength of schedule)
- Without using previous season's data (i.e. Regression).

Guess their end of season win percentage.

Naive guess (ask any rando on the street):

$$\widehat{WP} = \frac{W}{W+L}$$

What's wrong with this?

When Dodgers have only played a few games, this estimate is bad.

Ex $W=3, L=0, \widehat{WP}=1$

Idea Add fake data.

Suppose the Dodgers begin the season with W' wins and L' losses.

New guess:

$$\hat{WP}' = \frac{W+W'}{W+W'+L+L'}$$

For concreteness:

$$W=3, L=0, \quad \hat{WP}=1$$

$$W=3, L=0, W'=15, L'=15, \quad \hat{WP}' = \frac{18}{33} \approx .55$$

Quite different prediction early in the season

$$W=45, L=30, \quad \hat{WP} = \frac{45}{75} = .6$$

$$W=45, L=30, W'=15, L'=15, \quad \hat{WP}' = \frac{60}{105} \approx .57$$

similar prediction late in the season

Which is better?

Formalize this

Dodgers play $n=162$ games in a season.

Suppose, for simplicity, that the Dodgers win each game with probability p .

Game outcomes $\{x_1, \dots, x_n\}$, where

$$x_i \sim \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases} \stackrel{d}{=} \text{Bernoulli}(p)$$

Suppose we have observed m games thus far in the season.

Observed data $\{x_1, \dots, x_m\}$. Each x_i is 1 or 0.

Observed # wins $W = \sum_{i=1}^m x_i$.

So, $W \sim \text{Binomial}(m, p)$

$m = \# \text{ trials (games)}$
 $p = \text{prob. success (win)}$

and end-of-season win percentage $WP \sim \frac{1}{n} \text{Binomial}(n, p)$

Idea: use observed data to estimate p , call it \hat{p}

Then, estimate $\widehat{WP} = \frac{1}{n} \mathbb{E}[\text{Binomial}(n, \hat{p})] = \frac{1}{n} \cdot n\hat{p} = \hat{p}$.

Maximum Likelihood estimate (MLE)

choose \hat{p} to be the value of p which maximizes the probability of observing the game outcomes $\{x_1, \dots, x_m\}$ that we observed.

$$\hat{P}_{MLE} = \underset{p}{\operatorname{argmax}} \quad P(x_1, \dots, x_m \mid p)$$

likelihood : $P(\text{data given parameter})$

$$= \underset{p}{\operatorname{argmax}} \quad P(x_1 \mid p) \cdot P(x_2 \mid p) \cdot \dots \cdot P(x_m \mid p)$$

by independence

$$= \underset{p}{\operatorname{argmax}} \quad \prod_{i=1}^m P(x_i \mid p)$$

by def of product

$$= \underset{p}{\operatorname{argmax}} \quad \prod_{i=1}^m p^{x_i} (1-p)^{1-x_i}$$

because $x_i \sim \text{BER}(p)$

$$x_i=1 \text{ means } p^{x_i} (1-p)^{1-x_i} = p$$

$$x_i=0 \text{ means } p^{x_i} (1-p)^{1-x_i} = 1-p$$

$$= \operatorname{argmax}_p P^{\sum_{i=1}^m x_i} (1-p)^{\sum_{i=1}^m (1-x_i)}$$

$$= \operatorname{argmax}_p p^W (1-p)^L$$

where $W = \sum_{i=1}^m x_i = \text{number of wins (ones)}$
 $L = \sum_{i=1}^m (1-x_i) = \text{number of losses (zeros)}$

$$= \operatorname{argmax}_p \log [p^W \cdot (1-p)^L]$$

because \log is monotonic increasing
 to maximize $f(p)$ it to maximize $\log f(p)$

$$= \operatorname{argmax}_p W \log p + L \log (1-p)$$

to maximize the function $p \mapsto W \log p + L \log (1-p)$
 take the derivative and set it equal to 0
 (and check that the 2nd derivative is negative).

$$\frac{d}{dp} [W \log p + L \log (1-p)]$$

$$= W \cdot \frac{1}{P} - L \cdot \frac{1}{1-P} = 0$$

$$\Rightarrow \frac{W}{P} = \frac{L}{1-P} \Rightarrow P = \frac{W}{L}(1-P)$$

$$\Rightarrow P(1 + \frac{W}{L}) = \frac{W}{L} \Rightarrow P = \frac{\frac{W}{L}}{1 + \frac{W}{L}}$$

$$\Rightarrow \hat{P}_{MLE} = \frac{W}{W+L}$$

same formula
from earlier !!

The MLE is simply the observed win percentage midway through the season!

But we know this is a bad estimate early in the season.

So, why did the MLE go wrong??

How do we add the fake data W', L' to the MLE to get $\frac{W+W'}{W+W'+L+L'}$??

Before, to improve our estimate of WP, we added some fake number of wins W' and some fake number of losses L' to the Dodgers' record.

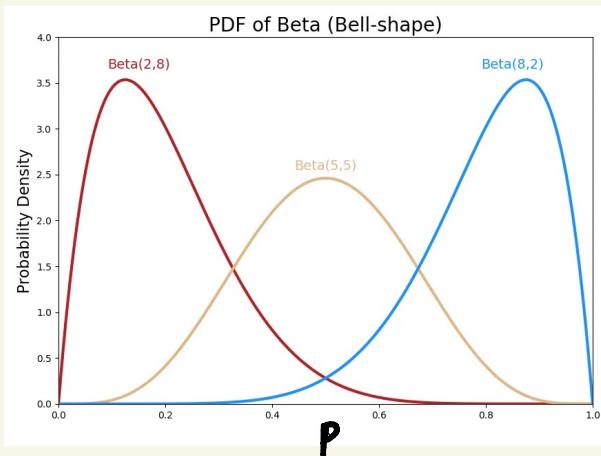
In other words, we assumed Prior Information: prior to the season, assume the Dodgers have W' wins and L' losses.

Formally, we use the Beta-Binomial model:

$$\begin{cases} W \sim \text{Binomial}(m, P) \\ P \sim \text{Beta}(\alpha, \beta) \\ \alpha = W', \quad \beta = L' \end{cases}$$

Beta distribution has density $f(p|\alpha, \beta) = C \cdot p^{\alpha-1} (1-p)^{\beta-1}$

on the interval $p \in [0, 1]$, where C is a constant chosen so that the distribution integrates to 1.



For example, $P \sim \text{Beta}(5,5)$ encodes a preference that P is closer to 0.5

As before, we wish to estimate P , this time with a Maximum a Posteriori (MAP) Estimate:

$$\hat{P}_{\text{MAP}} = \underset{P}{\operatorname{argmax}} \quad P(P|w)$$

Posterior = $P(\text{Parameter} | \text{data})$

$$= \underset{P}{\operatorname{argmax}} \quad \frac{P(w|P) \cdot P(P)}{P(w)} \quad \text{by Bayes' Rule}$$

$$= \underset{P}{\operatorname{argmax}} \quad \underbrace{P(w|P)}_{\text{likelihood}} \cdot \underbrace{P(P)}_{\text{prior}}$$

Since $P(w)$ has no P term

$$= \operatorname{argmax}_p P(\text{Binomial}(m, p) = w) \cdot P(\text{Beta}(\alpha, \beta) = p)$$

$$= \operatorname{argmax}_p \binom{m}{w} p^w (1-p)^{m-w} \cdot C p^{\alpha-1} (1-p)^{\beta-1}$$

$$= \operatorname{argmax}_p p^w (1-p)^L \cdot p^{\alpha-1} (1-p)^{\beta-1}$$

Since the constants $\binom{m}{w}$ and C don't have p terms

$$= \operatorname{argmax}_p p^{w+\alpha-1} (1-p)^{L+\beta-1}$$

••• same process as before

$$= \frac{w+\alpha-1}{w+\alpha-1 + L+\beta-1}$$

$$= \frac{w+w'}{w+w'+L+L'} \quad \text{if} \quad w' = \alpha-1 \\ L' = \beta-1$$

The MAP estimate is simply the win percentage if we add $\alpha-1$ fake wins and $\beta-1$ false losses!!

Note: $\alpha=1, \beta=1 \rightarrow \hat{P}_{\text{MAP}} = \hat{P}_{\text{MLE}}$
add no fake data

Can use past seasons to tune a smart choice for α, β .

HW: use past seasons' data to find good values of α and β , assuming $\alpha=\beta$ (i.e. assuming each team has a prior win percentage of 50%)

Tom Tango: $W=L=15$ is good

HW: use past seasons' data and pre-season win totals to tune smart values of α and β for each team

Takeaway: blending observed data

with prior knowledge is often a smart way to make better predictions which don't overfit to random noise