**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.

$$\underline{EX} \quad W = 3, \quad L = 0, \quad \widehat{WP} = 1$$

## Idea Add fake data.

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

New guess:

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

For concreteness:

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

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

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

quite different predictions early in the season

$W=45, L=30, \qquad \widehat{WP} = \frac{45}{75} = .6$

$W=45, L=30, W'=15, L'=15, \qquad \widehat{WP}' = \frac{60}{105} \approx .67$

similar predictions 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 $\{X_1, \ldots, X_n\}$, where

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

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

Observed data $\{X_1, \ldots, X_m\}$. Each $X_i$ is 1 or 0.

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

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

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

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

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

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

# Maximum Likelihood estimate (MLE)

Choose $\hat{p}$ to be the value of $p$ which maximizes the probability of observing the game outcomes $\{X_1, \ldots, X_m\}$ that we observed.

$$\hat{P}_{MLE} = \underset{p}{\arg\max} \; \mathbb{P}\left(X_1, \ldots, X_m \middle| P\right)$$

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

$$= \underset{p}{\arg\max} \; \mathbb{P}(x_1|P) \cdot \mathbb{P}(x_2|P) \cdot \ldots \cdot \mathbb{P}(x_m|P)$$

by independence

$$= \underset{p}{\arg\max} \; \prod_{i=1}^{m} \mathbb{P}(X_i|P)$$

by def of product

$$= \underset{p}{\arg\max} \; \prod_{i=1}^{m} p^{x_i} (1-p)^{1-x_i}$$

because $X_i \sim BER(P)$

$X_i = 1$ means $p^{x_i}(1-p)^{1-x_i} = P$

$X_i = 0$ means $p^{x_i}(1-p)^{1-x_i} = 1-P$

$$= \underset{p}{\text{argmax}} \quad p^{\sum_{i=1}^{m} x_i} \, (1-p)^{\sum_{i=1}^{m}(1-x_i)}$$

$$= \underset{p}{\text{argmax}} \quad p^{W} \, (1-p)^{L}$$

where $W = \sum_{i=1}^{m} x_i$ = number of wins (ones)

$L = \sum_{i=1}^{m} (1-x_i)$ = number of losses (zeros)

$$= \underset{p}{\text{argmax}} \quad \log \left[ p^{W} \cdot (1-p)^{L} \right]$$

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

$$= \underset{p}{\text{argmax}} \quad W \log p + L \log (1-p)$$

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).

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

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

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

$$\implies P\left(1 + \frac{W}{L}\right) = \frac{W}{L} \implies P = \frac{\frac{W}{L}}{1 + \frac{W}{L}}$$

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

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 data $(W', L')$.

In adding fake data, we used **prior information:** Prior to the season, we assumed the Phillies have $W'$ wins and $L'$ losses.

What is a way of formalizing prior information?

Bayesian statistics — the belief/philosophy that we should treat a parameter (e.g. $p$) as having a probability distribution

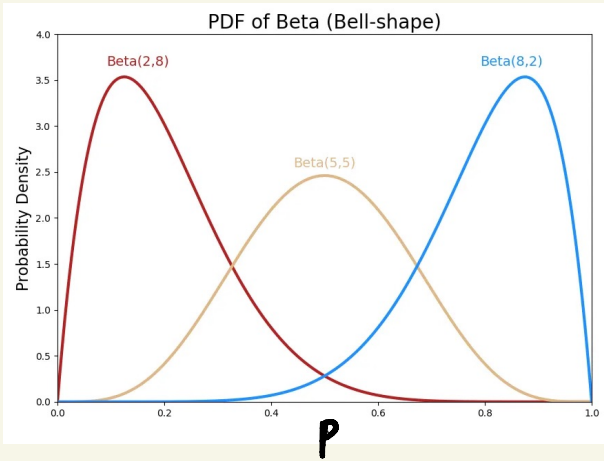Frequentist Statistics — treats a parameter as an unknown fixed number

So, our way of formalizing the addition of prior "fake" data is to, prior to seeing the data, give a probability distribution to the parameter (e.g. $p$) which reflects our prior belief on what $p$ is more likely to be than not.

Formally, we use the __Beta-Binomial model__ :

$$
\begin{cases}
W \sim Binomial(m, P) \\
P \sim Beta(\alpha, \beta) \quad \rightarrow PRIOR \\
\quad\quad \alpha = W' + 1, \quad \beta = L' + 1
\end{cases}
$$

__Beta distribution__ has density $f(p \mid \alpha, \beta) = C \cdot P^{\alpha - 1} (1-P)^{\beta - 1}$

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



PDF of Beta (Bell-shape)

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

As before, we wish to estimate $p$, this time with a
<u>Maximum a-Posteriori (MAP) Estimate:</u>

Choose the $\hat{p}$ which maximizes the Posterior
Probability of $p$.

Bayesian Approach
to Parameter
Estimation
$\left\{ \begin{array}{l} \text{1. PRIOR} \\ \text{2. observe data} \\ \text{3. adjust our posterior dist for} \\ \phantom{3.}\ p \text{ given the data} \end{array} \right.$

$$\hat{P}_{MAP} = \underset{p}{argmax}\ \underbrace{\mathbb{P}(p|w)}_{\text{Posterior} = \mathbb{P}(\text{parameter} | \text{data})}$$

$$= \underset{p}{argmax}\ \frac{\mathbb{P}(w|p) \cdot \mathbb{P}(p)}{\mathbb{P}(w)} \qquad \text{by Bayes' Rule}$$

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

since $\mathbb{P}(w)$ has no $p$ term

$$= \underset{p}{argmax}\ \mathbb{P}\left(\text{Binomial}(m,p) = w\right) \cdot \mathbb{P}\left(\text{beta}(\alpha,\beta) = p\right)$$

$$= \underset{p}{\text{argmax}} \quad \binom{m}{w} p^w (1-p)^{m-w} \cdot C \, p^{\alpha-1} (1-p)^{\beta-1}$$

$$= \underset{p}{\text{argmax}} \quad p^w (1-p)^L \cdot p^{\alpha-1} (1-p)^{\beta-1}$$

$$= \underset{p}{\text{argmax}} \quad p^{W+\alpha-1} (1-p)^{L+\beta-1}$$

$= \bullet\bullet\bullet$ same process as before

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

$$= \frac{W + W'}{W + W' + L + L'} \qquad \text{if} \quad \begin{aligned} W' &= \alpha - 1 \\ L' &= \beta - 1 \end{aligned}$$

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

Can use past seasons to tune a smart choice for $\alpha, \beta$.

Note: $\alpha = 1, \beta = 1 \implies \hat{P}_{MAP} = \hat{P}_{MLE}$

add no fake data

Model $\begin{cases} W \sim \text{Binomial}(m, p) \\ p \sim \text{Uniform}(0, 1) \end{cases} \longrightarrow$ uninformative PRIOR which encodes no preference on $p$

$$\hat{p}^{(MAP)} = \underset{p}{\text{argmax}} \; P(p \mid W) = \underset{p}{\text{argmax}} \; P(W \mid p) \cdot \underbrace{P(p)}_{1}$$

$$= \underset{p}{\text{argmax}} \; P(W \mid p) = \hat{p}^{(MLE)} = \frac{W}{W + L}$$

## Takeaways

- Bayesian Statistics: treat a parameter (e.g., $p$) as having a distribution

- Blend observed data with PRIOR knowledge, encoding info not seen in the data, to make better predictions