

The Power of Fake Data (PRIORS)

Q Suppose the Phillies have won W games and lost L games thus far in the season. How would you predict their end of season win percentage WP ?

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

Problem?
 - injuries } → ignore
 - no strength of schedule }
 - lack of data → $W=3, L=0, \widehat{WP}=1$

Idea Add Fake Data!

$$\widehat{WP}' = \frac{W+W'}{W+W'+L+L'} \quad \text{fake } (W', L')$$

Tom Tango: $W'=15, L'=15$

$$W=3, L=0, \widehat{WP}' = \frac{18}{33} \approx .55, \widehat{WP}=1$$

Some sort of regression happening here

Shrinking obs. w.p. $\frac{3}{3+0}$ to 50% $\frac{15}{15+15}$

Can we formalize this?

Phillies play $n=162$ games

Simple model: Phillies win each game w.p. p

Game outcomes X_1, \dots, X_n

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

We've observed just m of these game outcomes

X_1, \dots, X_m .

Observed # wins

$$W = \sum_{i=1}^m X_i \sim \text{Binomial}(m, p)$$

* At the same time, the actual end-of-season w.p. will be

$$\frac{1}{n} \text{Binomial}(n, p) = \frac{1}{n} \sum_{i=1}^n X_i$$

has expected value p .

So our task is to estimate p from our observed games.

* We will discuss a few ways to estimate p .

Maximum Likelihood Estimate (MLE)

Choose the \hat{p} which maximizes the probability of observing the data that we observed.

$$\hat{p}^{(MLE)} = \operatorname{argmax}_p \underbrace{P(x_1, \dots, x_m | p)}_{\text{Each } X_i \text{ Bernoulli}(p)}$$

$$= \operatorname{argmax}_p P(x_1 | p) \cdot P(x_2 | p) \cdot \dots \cdot P(x_m | p)$$

by independence

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

Product

$$P(x_i | p) \begin{cases} \rightarrow P(x_i=0 | p) = 1-p \\ \rightarrow P(x_i=1 | p) = p \end{cases}$$

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

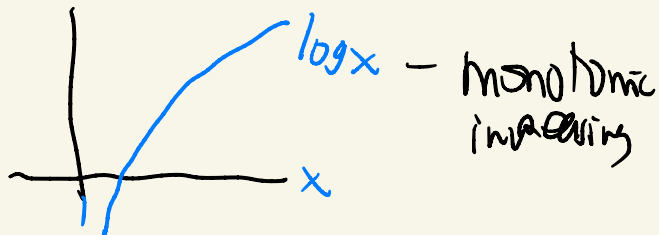
$$= \operatorname{argmax}_p \prod_{i=1}^m p^{x_i} (1-p)^{1-x_i}$$

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

$$\begin{aligned} \prod_{i=1}^m p^{x_i} (1-p)^{1-x_i} &= p^{x_1} p^{x_2} \dots p^{x_m} (1-p)^{1-x_1} (1-p)^{1-x_2} \dots (1-p)^{1-x_m} \\ &= p^{x_1+x_2+\dots+x_m} (1-p)^{(1-x_1)+(1-x_2)+\dots+(1-x_m)} \end{aligned}$$

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

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



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

How do we maximize this?
set derivative equal to 0 and solve.

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

$$\Rightarrow W \frac{1}{p} = L \frac{1}{1-p}$$

$$W(1-p) = Lp$$

$$W = p(W+L)$$

$$p = \frac{W}{W+L}$$

$$\hat{p}^{(MLE)} = \operatorname{argmax}_p P(x_1, \dots, x_m | p) = \frac{W}{W+L}$$

We know this is bad.

Why did the MLE go wrong?

Vaguely though we know we are looking to add **fake data** (w', l') to achieve
$$\frac{w+w'}{w+w'+l+l'}$$
.

In adding fake data, we use

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 a parameter (e.g. p) itself has an (unknown) distribution

Frequentist Statistics — treats a parameter
as a fixed (unknown) number

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

$$p \sim$$

Prior
Distribution

Our way of formalizing the addition of prior "fake data" is to, prior to seeing any 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.

Targ: added $W' = 15, L' = 15 \rightarrow p \approx \frac{15}{30} = \frac{1}{2}$

$$\hat{w}_p = \frac{W + 15}{W + L + 30}$$

early in sym closer to $1/2$
later in sym closer to $\text{dr. } W$

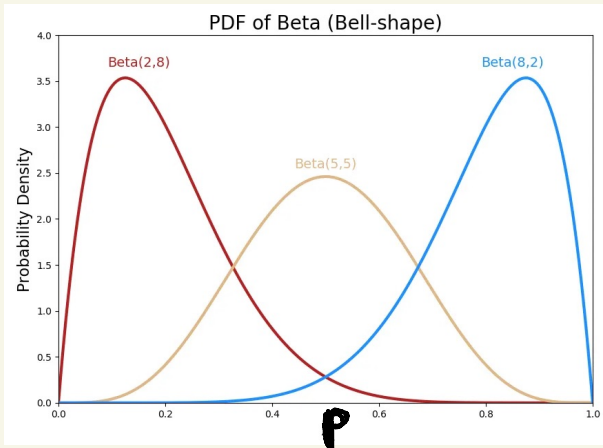
Prior: loosely, lets be closer to $\frac{1}{2}$ and further from the extremes 0 and 1

What prior dist. should we use for p ?
 $p \in [0, 1]$

Simplest dist. defined on $[0, 1]$
Which has flexible parameters? **Beta**

$$\begin{cases} W \sim \text{Binomial}(n, p) \\ P \sim \text{Beta}(\alpha, \beta) \end{cases} \rightarrow \text{Prior Distribution}$$

Beta-Binomial Model



Y-axis: $P(\text{Beta}(\alpha, \beta) = p)$

$p \sim \text{Beta}(\alpha, \beta)$

$p \in [0, 1]$

Beta-dist. has density

$$P(\text{Beta}(\alpha, \beta) \in [p, p+dp]) = f(p|\alpha, \beta)$$

$$= C \cdot p^{\alpha-1} \cdot (1-p)^{\beta-1}$$

where C is a constant

so that $\int_0^1 f(p) dp = 1.$

Have Model $\begin{cases} W \sim \text{Binomial}(n, p) \\ P \sim \text{Beta}(\alpha, \beta) \end{cases} \rightarrow \text{Prior Distribution}$

Want to estimate $p.$

Before we ever introduced a prior:

Maximum Likelihood Estimate (MLE)

Choose the \hat{p} which maximizes the probability of observing the data that we observed.

$$\hat{p}^{(MLE)} = \underset{p}{\text{argmax}} P(x_1, \dots, x_n | p)$$

Now that have a prior:

Maximum a-Posteriori (MAP)

Choose the \hat{p} which maximizes the posterior probability of p (the probability of p using the information from our observed data).

Bayesian Approach to parameter estimation

1. prior $p \sim \text{Beta}(\alpha, \beta)$

2. observe data X_1, \dots, X_m

3. adjust our dist. for p using the data

$$P(p | X_1, \dots, X_m)$$

$$\hat{p}^{(\text{MAP})} = \underset{p}{\text{argmax}} \underbrace{P(p | W, L)}$$

the posterior prob. of a parameter is the prob. of that parameter given the observed data

Bayes Rule

$$= \operatorname{argmax}_p \frac{P(W, L | p) \cdot P(p)}{P(W, L)}$$

$$= \operatorname{argmax}_p \underbrace{P(W, L | p)}_{\text{Binomial}} \cdot \underbrace{P(p)}_{\text{Beta}}$$

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

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

$$P(W, L | p) = P\left(\begin{array}{l} \text{obs.} \\ \# \text{ wins} = W, \\ \text{obs.} \\ \# \text{ losses} = L \end{array} \middle| p\right)$$

$$= P\left(\sum_{i=1}^m X_i = W, m - \sum_{i=1}^m X_i = L \middle| p\right)$$

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

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

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

$$\begin{aligned} \alpha - 1 &= W' \\ \beta - 1 &= L' \end{aligned}$$

$$\left\{ W \sim \text{Binomial}(m, p) \right\} \Rightarrow \hat{p}^{(\text{MLE})} = \frac{W}{W+L}$$

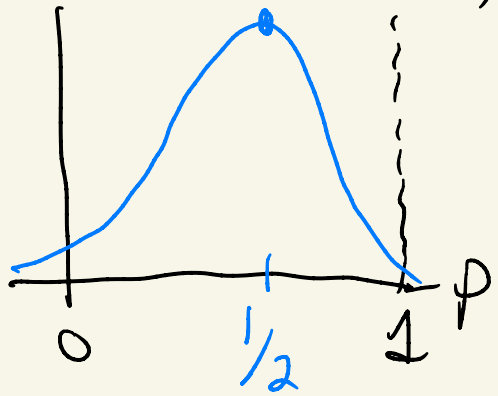
$$\left\{ \begin{array}{l} W \sim \text{Binomial}(m, p) \\ p \sim \text{Beta}(\alpha, \beta) \\ \alpha = W' + 1, \beta = L' + 1 \end{array} \right\} \Rightarrow \hat{p}^{(\text{MAP})} = \frac{W + W'}{W + W' + L + L'}$$

$$\left\{ \begin{array}{l} W \sim \text{Binomial}(m, p) \\ p \sim \text{Uniform}(0, 1) \end{array} \right\} \Rightarrow \hat{p}^{(\text{MAP})} = \hat{p}^{(\text{MLE})} = \frac{W}{W+L}$$

Takeaways

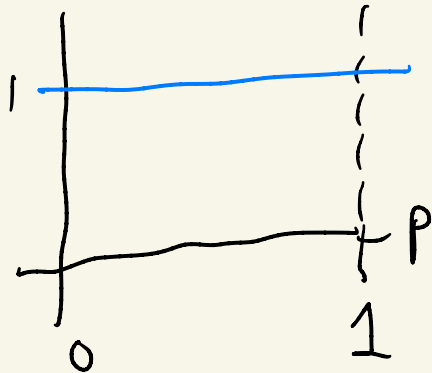
- Bayesian approach: treat parameter (e.g. p) as having a distribution
- PRIOR: allows us to make better predictions by encoding information not seen in the available data

* If you have prior reason to believe, before looking at data, that p should be $\approx \frac{1}{2}$ when small # of obs., then choose a prior like $\text{Beta}(16, 16)$ which looks like



* Conversely, if you have no prior info, or beliefs on whether p should be,

$\text{Uniform}[0, 1]$



model $\begin{cases} W \sim \text{Binomial}(n, p) \\ p \sim \text{Uniform}[0, 1] \end{cases}$ prior

$$\begin{aligned} \hat{p}^{(\text{MAP})} &= \underset{p}{\text{argmax}} P(p | W) \\ &= \underset{p}{\text{argmax}} P(W | p) \cdot \underbrace{P(p)} \end{aligned}$$

$$P(p) = P(\text{Uniform}[0, 1] = p) = 1$$

$$= \underset{p}{\text{argmax}} P(W | p)$$

$$= \hat{p}^{(\text{MLE})} = \frac{W}{W+L}$$