

# lec 10 MATH 341/641

We say the Beta distr is the "conjugate prior" for the iid Bernoulli OGP. Conjugacy means prior and posterior are same rv (but different parameter values updated by the data).

$$\text{Beta}(\alpha, \beta) \xrightarrow{\vec{x}} \text{Beta}(\underbrace{\alpha + \sum x_i}_{\# \text{ successes}}, \underbrace{\beta + n - \sum x_i}_{\# \text{ failures}})$$

# pseudosuccesses      # pseudofailures

What is the interpretation of the values of the hyperparameters?

The prior's hyperparameters are like observing fake data,  $n_0 = \alpha + \beta$  <sup>where</sup> # pseudotrials

If we employ Laplace's prior of indifference  $\mathcal{U}(0,1) = \text{Beta}(\alpha=1, \beta=1)$

$\Rightarrow n_0 = 2$  pseudotrials. This is weird. Further, consider

$$\begin{aligned} \hat{\theta}_{\text{MSE}} &= \frac{\alpha + \sum x_i}{\alpha + \beta + n} = \frac{\alpha}{\alpha + \beta + n} \cdot \frac{\alpha + \beta}{\alpha + \beta} + \frac{\sum x_i}{\alpha + \beta + n} \cdot \frac{n}{n} \\ &= \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} + \frac{n}{\alpha + \beta + n} \frac{\sum x_i}{n} \\ &= \rho \underbrace{E[\theta]}_{\text{prior guess}} + (1 - \rho) \underbrace{\hat{\theta}_{\text{MLE}}}_{\text{pure data estimate}} \end{aligned}$$

$$\lim_{n \rightarrow \infty} \hat{\theta}_{\text{MSE}} = \hat{\theta}_{\text{MLE}} \quad \text{since} \quad \lim_{n \rightarrow \infty} \rho = 0$$

$$\text{let } \rho := \frac{\alpha + \beta}{\alpha + \beta + n},$$

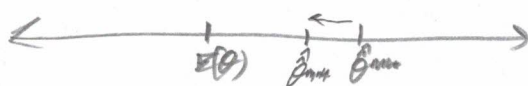
the shrinkage metric, measures strength of prior  $\in [0,1]$  on the most popular point estimate.

e.g.  $\alpha = \beta = 1, n = 3 \Rightarrow \rho = \frac{2}{5}$

$\Rightarrow$  40% weight on prior, 60% weight on data

$\hat{\theta}_{\text{MSE}}$  is called a "shrinkage estimator"

since it shrinks toward  $E(\theta)$



13

Prob:  $X \sim \text{bin}(n, \theta)$ ,  $f(\theta) = \text{Beta}(\alpha, \beta)$

$$\begin{aligned} f(\theta|\bar{x}) &\propto P(X|\theta) f(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \mathbb{I}_{\theta \in (0,1)} \\ &\propto \theta^{\alpha+x-1} (1-\theta)^{\beta+n-x-1} \mathbb{I}_{\theta \in (0,1)} \\ &\propto \text{Beta}(\alpha+x, \beta+n-x) \end{aligned}$$

Same as before although  $x = \sum x_i$  which is exactly what it is

Since  $\text{Binomial}(n, \theta) = \text{sum of } n \text{ iid Bern}(\theta)$

(3)

Let's now ask a different question: You see the  $n$  trials and now you want to predict how the future  $n_x$  trials are distributed.

$$\frac{0}{1} \quad \frac{1}{2} \quad \dots \quad \frac{1}{n} \quad \bigg| \quad \frac{1}{1} \quad \frac{1}{2} \quad \dots \quad \frac{1}{n_x}$$

How is  $X_x$  distributed?

This is called "prediction" or "forecasting" and is mainly the subject of MATH 342/642. Statistical Inference focuses on explanation of models, not prediction. But Bayesian inference gives you it for free!

Obviously if  $\theta$  is known,  $X_x \sim \text{Bin}(n_x, \theta)$ .

So what's the best guess  $X_x \sim \text{Bin}(n_x, \hat{\theta}_{MLE})$ ?

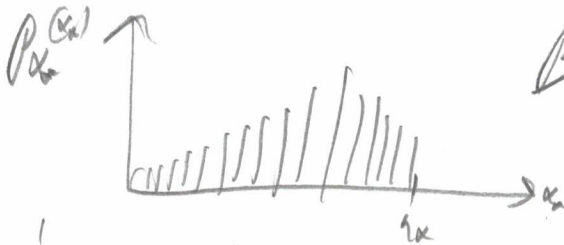
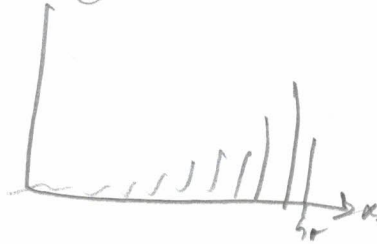
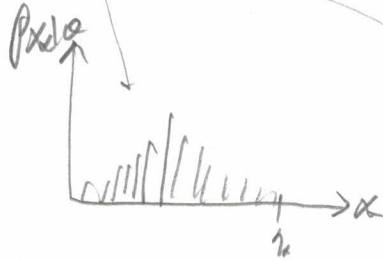
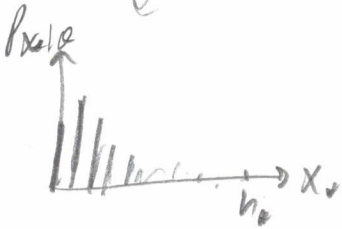
Yes, in frequentist statistics, this is the best you can do, but it ignores the uncertainty in the estimate  $\hat{\theta}_{MLE}$ ! So it's not really correct.

In Bayesian, you have an idea of the uncertainty in  $\theta$  after the data is seen. It's called the posterior!

$$P(\theta|x) \propto \text{Beta}(\alpha+x, \beta+n-x)$$

4

Why not copy and paste this derivation  
using the OBP? Draw a  
prior  $\theta|x$ , then sample a  
Binomial( $n_0, \theta$ ). Then another,  
then another. This averages  
over the uncertainty in the  
knowledge of  $\theta$ .



BetaBinomial!

$$\frac{\binom{n}{x}}{\binom{n}{x_0}} \frac{B(\alpha+x, \beta+n-x)}{B(\alpha+x_0, \beta+n-x_0)}$$

$$P(x_0|x) = \int_0^1 \underbrace{P(\theta|x)}_{\text{Bin}(n, \theta)} \underbrace{f(\theta|x)}_{\text{Beta}(\alpha+x, \beta+n-x)} d\theta = \text{BetaBinomial}(n, \alpha+x, \beta+n-x)$$

Proof in MATH 340

Concrete Example. Baseball player has  $n=10$  at bats.  $x=6$  are hits.

What's the prob he will have  $x=17$  hits at the next  $n_0=37$  at bats?  
Assume Laplace prior

$$P(x_0|x=6) = \text{BetaBinomial}(37, 1+6, 1+4)$$

$$P(x_0=17|x=6) = \frac{\binom{37}{17}}{\binom{7}{5}} B(24, 25) = \text{dbernbinomial}(17, 37, 7.5) =$$

↑  
leaked for test

What's the prob he gets 17 hits or less on the test 37 at bats? 5

$$P(X_n \leq 17, n=37) = \text{P}(\text{binomial}(17, 37, 7, 5))$$

$\uparrow \quad \uparrow \quad \uparrow \quad \nwarrow$   
 $n \quad x \quad \alpha \quad \beta$

will realize to either value of

What if  $n=1$ ?  $X_n \sim \text{Bernoulli}(\theta_n)$  since it can be 0 or 1. What is  $\theta_n$ ?

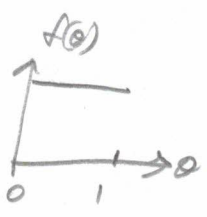
$$P(X_n | \alpha) = \text{BetaBinom}(1, \alpha + x, \beta + n - x) = \frac{\binom{1}{x}}{\text{Beta}(\alpha + x, \beta + n - x)} \text{Beta}(\alpha + x, 1 + x + \beta + n - x)$$

$$\begin{aligned} \theta_n &= P(X_n = 1 | \alpha) = \frac{\binom{1}{1}}{\text{Beta}(\alpha + 1, \beta + n - 1)} \text{Beta}(1 + \alpha + 1, \beta + n - 1) \\ &= \frac{1}{\frac{\Gamma(\alpha + 1) \Gamma(\beta + n - 1)}{\Gamma(\alpha + \beta + n)}} \cdot \frac{\Gamma(1 + \alpha + 1) \Gamma(\beta + n - 1)}{\Gamma(1 + \alpha + \beta + n)} = \frac{\alpha + 1}{\alpha + \beta + n} = \hat{\theta}_{\text{MLSE}} \end{aligned}$$

The prob the test <sup>at</sup> bat we will get hit is  $\frac{\alpha + 1}{\alpha + \beta + n}$

Priors are subjective. But can we make them "dejective" to minimize their effect on the inference. Have cake and eat it too.

Laplace's prior is "flat" i.e. it doesn't give any special preference to any of the  $\theta$ 's.



But we also say it implies  $E(\theta) = 0.5$  and  $v_0 = 2$  i.e.  $x < 1$ , or  $x = 0$ .

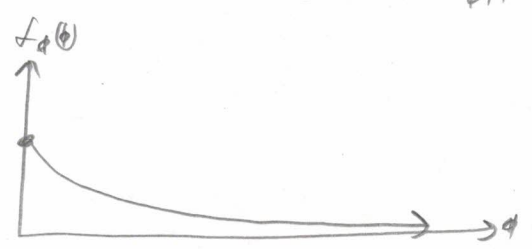
$$f(\theta|x) = \text{Beta}(\alpha + x, b + n - x)$$

That doesn't seem too dejective! Further...

Let's say we care about inference for  $\phi = g(\theta) = \frac{1-\theta}{\theta}$ , odds against who does Laplace's prior mean on this new scale? 1:1 transformation

$$f_{\phi}(\phi) = f_{\theta}(g^{-1}(\phi)) \left| \frac{d}{d\phi} [g^{-1}(\phi)] \right| = \frac{1}{\phi+1} \mathbb{1}_{\phi+1 \in [0,1]} \frac{1}{(\phi+1)^2} = \frac{1}{(\phi+1)^2} \mathbb{1}_{\phi+1 \in [1,\infty)} = \frac{1}{(\phi+1)^2} \mathbb{1}_{\phi \in [0,\infty)}$$

$$\phi = \frac{1}{\theta} - 1 \Rightarrow \phi + 1 = \frac{1}{\theta} \Rightarrow \theta = \frac{1}{\phi+1} = g^{-1}(\phi) \left| \frac{d}{d\phi} \left[ -(\phi+1)^{-2} \right] \right| = \frac{1}{(\phi+1)^2} \propto F_{2,2} \text{ dist} = \text{BetaPrime}(1, 1)$$



Which is not flat!! Further, there is no way to have a flat PDF on  $(0, \infty)$ ! Inference prior Laplace is fixed to a specific parameterization

You are saying small odds against values are more likely than large

What happens in this case??