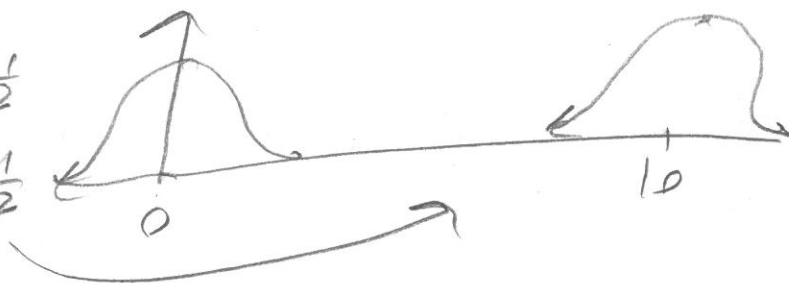


Lecture 7 2/24/16

11

If

$$X \sim \begin{cases} N(0, 1^2) & \text{w.p. } \frac{1}{2} \\ N(10, 1^2) & \text{w.p. } \frac{1}{2} \end{cases}$$



What is PDF of  $X$ ?

$$= \frac{1}{2} N(0, 1) + \frac{1}{2} N(10, 1^2)$$

we will write this later

$$p(x) = \sum_{\theta \in \{0, 10\}} p(x; \theta) p(\theta) = \sum_{\theta \in \{0, 10\}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\theta)^2} \left(\frac{1}{2}\right)$$

$$= \frac{1}{2\sqrt{2\pi}} \left( e^{-\frac{1}{2}x^2} + e^{-\frac{1}{2}(x-10)^2} \right)$$

Okay... What about

$$X \sim \begin{cases} \text{Bin}(n, 0.25) & \text{w.p. } 0.1 \\ \text{Bin}(n, 0.4) & \text{w.p. } 0.9 \end{cases}$$

$$p(x) = \sum_{\theta \in \{0.25, 0.4\}} \binom{n}{x} \theta^x (1-\theta)^{n-x} p(\theta)$$

"marginal distribution"  
or  
"compound distr."

$$= \binom{n}{x} \left( .25^x .75^{n-x} \cdot 0.1 + .4^x .6^{n-x} \cdot 0.9 \right)$$

What if

$$X \sim \text{Bin}(n, \theta) \text{ but } \theta \sim U(0, 1)$$

$$\begin{pmatrix} RRR \\ GGG \end{pmatrix} \quad \begin{pmatrix} RGG \\ GGG \end{pmatrix} \quad \begin{pmatrix} RRR \\ RRG \end{pmatrix} \dots$$

Bags of marbles each with  
diff prop of Red's  
the bags form a  $U(0, 1)$

$$p(x) = \int_{\theta \in (0,1)} p(x; \theta) p(\theta) d\theta = \int_0^1 \binom{n}{x} \theta^x (1-\theta)^{n-x} d\theta = \binom{n}{x} B(x+1, n-x+1)$$

this is a dist!

Wait what does this  
look like?

denominator is by

$$\sum_{x=0}^n \binom{n}{x} B(x+1, n-x+1) = 1$$

→ will see who is called soon

What about  $X \sim \text{Bin}(n, \theta)$  but  $\theta \sim \text{Beta}(\alpha, \beta)$ ?

$$p(x) = \int_0^1 \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta = \frac{\binom{n}{x}}{B(\alpha, \beta)} \int_0^1 \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta$$

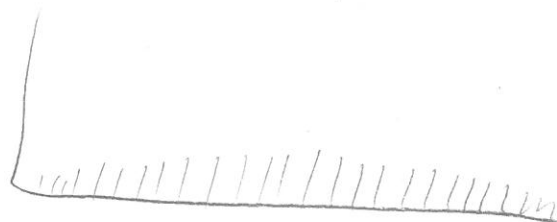
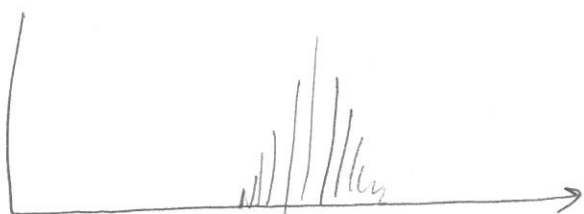
$$= \binom{n}{x} \frac{B(x+\alpha, n-x+\beta)}{B(\alpha, \beta)}$$

"BetaBin( $n, \alpha, \beta$ )"

[3]

Known as the beta-binomial distribution.  $E(X)$ ?  $\frac{\alpha}{\alpha+\beta}$   $Var(X) = \frac{n\alpha\beta(\alpha+\beta+1)}{(\alpha+\beta)^2(\alpha+\beta+1)}$

$X \sim Bin(n, \theta)$  vs  $X \sim BetaBin(n, \alpha, \beta)$  Assume  $n$  is the same.



Bin has one param

$\Rightarrow E(X) = n\theta$   
 $Var(X) = n\theta(1-\theta)$  } pick mean  $\Rightarrow$  var fix

BetaBin has two param  $\Rightarrow$  you can pick mean & variance. You can get "burstiness" or "overdispersion" as above or very low variance

If  $\alpha, \beta \rightarrow \infty$  it can use

$$\lim E(X) = n \left( \frac{1}{2} \right)$$

$$\lim Var(X) = n \lim_{\alpha \rightarrow \infty} \frac{\alpha^2 (2\alpha+1)}{(2\alpha)^2 (2\alpha+1)} = n \frac{1}{4}$$

$$\Rightarrow \lim_{\alpha, \beta \rightarrow \infty} BetaBin(n, \alpha, \beta) = Bin(n, \frac{1}{2})$$

That makes sense because

likewise  $\lim_{\alpha, \beta \rightarrow \infty} \frac{\alpha}{\alpha+\beta} = c$  you can get any  $Bin(n, \theta)$

that just is  $Bin(n, \theta)$  where  $\theta$  is drawn from degenerate distr.

The overdispersion gives you flexibility in the case where you're unsure about  $\theta$ .

6/15

4

e.g. Wikipedia. Families with 13 children. Gender of first 12

6/15

Children	# of males	0	1	2	3	4	5	6	7	8	9	10	11	12	13
X		3	29	104	286	620	1033	1243	1112	829	428	181	45	7	
$\hat{X}$ Bern Bin(13, 0.5)		2	23	105	311	656	1036	1258	1182	854	462	178	44	5	
$\hat{X} \sim \text{Bin}(n, .5n)$		1	12	72	253	628	1585	1867	1266	854	410	132	26	2	

36,100 kids  
-519

bern bin fits better! Why is human gender ratio bern-bin on families??

Recall:  $X_1, \dots, X_n \overset{\text{exch}}{\sim} \text{Bern}(\theta), \theta \sim \text{Bern}(\alpha, \beta)$

$$\Rightarrow \theta | X \sim \text{Bern}(\alpha + x, \beta + n - x) \quad \& \quad X^* | X \sim \text{Bern}\left(\frac{x + \alpha}{n + \alpha + \beta}\right)$$

improved  $X \sim \text{BernBin}(n, \alpha, \beta) \Rightarrow$  prior predictive distr for dataset of size  $n$

$$X^* | X \sim ???$$

let  $X^*$  be length  $m$ .

We are predicting the outcome not of only the very next trial... but of the next  $m$  trials given  $n$  data pts and a bern prior or  $\theta$ .

$$P(X^* | x) = \int \underbrace{P(X^* | \theta)}_{\text{Bin}(m, \theta)} \underbrace{P(\theta | x)}_{\text{Bern}(\alpha + x, \beta + n - x)} d\theta$$

$\theta \in (0, 1)$

$$= \int_0^1 \binom{m}{x^*} \theta^{x^*} (1-\theta)^{m-x^*} \frac{1}{B(\alpha+x, \beta+x)} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1} d\theta$$

$$= \frac{\binom{m}{x^*}}{B(\alpha+x, \beta+x)} \int_0^1 \theta^{x^*+x+\alpha-1} (1-\theta)^{m-x^*+n-x+\beta-1} d\theta$$

$$= \frac{\binom{m}{x^*}}{B(\alpha+x, \beta+x)} B(x^*+x+\alpha, m-x^*+n-x+\beta)$$

$$= \text{Bern Bin}(m, x+\alpha, n-x+\beta) \quad \text{Hw: prove for } m=1$$

this is a beta

Why should this be?

If  $\theta$  was <sup>precisely</sup> known

$$X^* | \theta = X^* | \theta \sim \text{Bin}(m, \theta)$$

But  $\theta$  is known with uncertainty and the uncertainty  $\sim \text{Beta}$ ,  
it is as if on each trial, you grab a  $\theta | X$  and use this  
to run a Bernoulli, then grab another one, etc...

FACT:

Posterior predictive distribution / function / incorporation  
Uncertainty in your best guess of  $\theta$ !!

Review...  $X_1, \dots, X_n$  i.i.d.  $\text{Bern}(\theta)$ , or  $\text{Beta}(\alpha, \beta)$

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)} \propto P(x|\theta) P(\theta)$$

↑  
misses the  $\frac{1}{P(x)}$  "normalization factor"

Note that it's now a function of  $\theta$

$$P(\theta|x) \propto \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} = \frac{\binom{n}{x}}{B(\alpha, \beta)} \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

" $P(\theta|x; \alpha, \beta)$ "

↑↑↑

Assuming  
fixed

"hyperparameters"

"prior"

$$\propto \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

Why?

$\forall \theta$  this is a constant

$\forall \theta$  this is a constant. But what about  $x$ ?  $P(\theta|x)$ .  $x$  is fixed!!

That's what conditional prob means!

Compound distribution marginalizing

over the intermediate parameter  $\theta$

Learn the model as a prob of hyperparameters

$$P(x; \alpha, \beta) = \int P(x|\theta) P(\theta; \alpha, \beta) d\theta$$

$\theta \in (0,1)$

$$P(\theta|x) \propto \theta^{x+\alpha-1} (1-\theta)^{n-x+\beta-1}$$

$$\Rightarrow \frac{P(\theta_1|x)}{P(\theta_2|x)} = \left( \frac{\theta_1}{\theta_2} \right)^{x+\alpha-1} \left( \frac{1-\theta_1}{1-\theta_2} \right)^{n-x+\beta-1}$$

You can get posterior odds directly!!

But can you get the whole density? Yes...

7

$$P(\theta|x) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

If I integrate this... what do I get?

$$\int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta \neq 1 \text{ but } < \infty$$
$$\Rightarrow = C$$

$$\int g(x;\theta) d\theta = C \Rightarrow f(\theta) = \frac{g(x;\theta)}{C} \text{ so it is now a density}$$

Hold on...  $X \sim \text{Bin}(n, \theta)$

$$P(X|n) = \binom{n}{x} \theta^x (1-\theta)^{n-x} = \frac{n!}{(n-x)! x!} \theta^x (1-\theta)^n (1-\theta)^{-x}$$

$$\propto \frac{1}{(n-x)! x!} \left(\frac{\theta}{1-\theta}\right)^x$$

this doesn't look like a binomial!!

Kernel of the PMF/PDF

$$f(x; \theta) = \underbrace{h(\theta)}_{\text{normalization constant}} \underbrace{g(x, \theta)}_{\text{"Kernel"}}$$

normalization  
constant

Given a kernel, the normalization constant is fixed -  
Therefore a r.v. can be identified by its kernel!

$$\alpha, \beta \in (-1, \infty)$$

(2)

$$\begin{aligned} \mathcal{O} \sim \text{Beta}(\alpha, \beta) &:= \frac{1}{\Gamma(\alpha, \beta)} \mathcal{O}^{\alpha-1} (1-\mathcal{O})^{\beta-1} \propto \mathcal{O}^{\alpha-1} (1-\mathcal{O})^{\beta-1} \\ &= \mathcal{O}^a (1-\mathcal{O})^b \quad a, b \in (-1, \infty) \end{aligned}$$

Anything of this form is a beta!

$$p(\mathcal{O}|x) \propto \underbrace{\mathcal{O}^{\overbrace{x+\alpha-1}^a} (1-\mathcal{O})^{\overbrace{n-x+\beta-1}^b}}_{\text{the kernel for a beta distr.}}$$

$$\Rightarrow \mathcal{O}|x \sim \text{Beta}(a, b)$$

How: you will find the kernels for all the common distr's.

Deys-Lythe prior

$$\mathcal{O}|x \sim \text{Beta}(\alpha+x, \beta+n-x), \quad \mathcal{O} \sim U(0,1) = \text{Beta}(1,1)$$

"Wilson Estimate" /

$$\hat{\mathcal{O}}_{\text{mmsE}} = \frac{\alpha+x}{\alpha+\beta+n} = \frac{x+1}{n+2} \quad \left. \vphantom{\hat{\mathcal{O}}_{\text{mmsE}}} \right\} \text{"Law of Succession". Very famous}$$

estimate of  $\mathcal{O}$ . Highly biased only because Bayesian stat is highly biased or

$$\hat{\mathcal{O}}_{\text{map}} = \frac{x}{n} = \hat{\mathcal{O}}_{\text{MLE}} \quad (\text{interesting})$$

He was working on the Schrödinger problem posed by Hume.



Recall the post. distr:

(9)

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

$$\hat{\theta}_{\text{MMSE}} = E[\theta | x] = \frac{\alpha+x}{\alpha+\beta+n} = \frac{\alpha}{\alpha+\beta+n} + \frac{x}{\alpha+\beta+n}$$

$$\frac{\alpha}{\alpha+\beta+n} \left( \frac{\alpha+\beta}{\alpha+\beta} \right) + \frac{x}{\alpha+\beta+n} \left( \frac{n}{n} \right)$$

$$= \frac{\alpha+\beta}{\alpha+\beta+n} \left( \frac{\alpha}{\alpha+\beta} \right) + \frac{n}{\alpha+\beta+n} \left( \frac{x}{n} \right)$$

$$\hat{\theta}_{\text{MLE}} = \bar{x}$$

$$\text{Note } \frac{\alpha+\beta}{\alpha+\beta+n} + \frac{n}{\alpha+\beta+n} = 1$$

Let  $\rho := \frac{\alpha+\beta}{\alpha+\beta+n} \Rightarrow 1-\rho$  and both are prop's which total to 1.

$$\Rightarrow \hat{\theta}_{\text{MMSE}} = \rho E[\theta] + (1-\rho) \hat{\theta}_{\text{MLE}}$$

Where  $\rho$  is known as the shrinkage factor / proportion.

At  $\hat{\theta}_{\text{MMSE}}$  is known as a "shrinkage estimator"

Since it "shrinks" towards your prior best guess,  $E[\theta]$ .

Large values of  $\rho$  shrink more toward  $\rho$  and small values of  $\rho$

let the data speak for itself. If  $n \uparrow$ ,  $\rho \downarrow \Rightarrow$  data speaks for itself

$$\text{Beta}(\alpha+x, \beta+n-x) \propto \theta^{(\overset{\text{\# prior success}}{\downarrow} (\alpha-1) + \overset{\text{\# success}}{\downarrow} x) (1-\theta)^{(\overset{\text{\# prior failure}}{\downarrow} (\beta-1) + \overset{\text{\# failure}}{\downarrow} (n-x))} \quad (\text{How 3 \# 3c})$$

So  $\text{Dir}(\theta; 1) = \text{Beta}(1,1)$  is equivalent to seeing no success or failure

but  $\rho = \frac{2}{2+n} \rightarrow 0$  as  $n \rightarrow \infty$  but it's still shrinks toward  $E(\theta) = \frac{1}{2}$



10

