

Lec 6 2/20/18 Ruth 391

Previously

$\theta \sim p(\theta)$ the prior

X_1, X_2, X_3 observed

we want to know the distribution of X_4 given our (a) prior, (b) data

$$P(X_4 | X_1, X_2, X_3) = \sum_{\theta \in \Theta_0} P(X_4, \theta | X_1, X_2, X_3)$$

$$P(A \cap B) = P(A|B) P(B)$$

$$P(A, B | C) = P(A|B, C) P(B|C)$$

joint = cond x marginal

Looks messy but it's the same as

$$P(Y) = \sum_{X \in \text{supp}(X)} P(X, Y) \quad \text{margining out } X$$

$$= \sum_{\theta \in \Theta_0} P(X_4 | \theta, X_1, X_2, X_3) P(\theta | X_1, X_2, X_3)$$

" $P(X_4 | \theta)$ why?"

posterior

Since $X_1 | \theta, X_2 | \theta, X_3 | \theta, X_4 | \theta \stackrel{\text{iid}}{\sim} \text{Bern}(\theta)$

Pf:

def. Cond prob.

def of cond. prob.

$$P(X_4 | \theta, X_1, X_2, X_3) = \frac{P(\theta, X_4, X_1, X_2, X_3)}{P(\theta, X_1, X_2, X_3)} = \frac{P(X_4, X_1, X_2, X_3 | \theta) P(\theta)}{P(X_1, X_2, X_3 | \theta) P(\theta)} = \frac{P(X_4 | \theta) P(X_1 | \theta) P(X_2 | \theta) P(X_3 | \theta)}{P(X_1 | \theta) P(X_2 | \theta) P(X_3 | \theta)} = P(X_4 | \theta)$$

$$P(X_4 | X_1, X_2, X_3) = \sum_{\theta \in \Theta_0} P(X_4 | \theta) P(\theta | X_1, X_2, X_3)$$

How to parse this?

a weight

this is the same as $X_4 | \theta \sim \text{Bern}(\theta)$ but θ is averaged over all θ 's in the posterior, weighted by the posterior.

$\neq P(X_4 | \hat{\theta}_{MLE})$ is the best done with the frequentist's MLE. $X_1=0, X_2=0, X_3=0? \hat{\theta}_{MLE}=0 \Rightarrow X_4 \sim \text{Bern}(0)$ BAD!

Recall the old Bernoulli model

with $\Theta_0 = \{0.25, 0.75\}$ $P(\theta) = U(\Theta_0) = \begin{cases} 0.25 & \text{up } 0.5 \\ 0.75 & \text{up } 0.5 \end{cases}$

principle of indifference

since $P(\theta)$ const. under princ. of indifference

DATA: $X_1=0, X_2=1, X_3=1$

What is Bayesian estimate of θ ?

$$\hat{\theta}_{MAP} = \underset{\theta \in \Theta_0}{\operatorname{argmax}} \{P(\theta|x)\} = \underset{\theta \in \Theta_0}{\operatorname{argmax}} \{P(x|\theta)P(\theta)\} = \underset{\theta \in \Theta_0}{\operatorname{argmax}} \{P(x|\theta)\} = 0.75$$

most probable θ after seeing data in the support of the prior, Θ_0

but $0.75 \neq 0.66 = \hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{L(\theta; x)\} = \underset{\theta \in \Theta}{\operatorname{argmax}} \{P(x; \theta)\}$

Why not? $\Theta_0 \neq \Theta$

likely a bad idea - why should you put zero prob on pieces of the parameter space?

So let's say $\Theta_0 = \Theta = (0,1)$, the parameter space of the Bernoulli

What prior should we put on this? Via the principle of indifference, we use the

$\theta \sim U(0,1)$ i.e. the 'uniform prior'. It is non continuous
stochastic uniform

Recall its density $P(\theta) = \begin{cases} 1 & \text{if } \theta \in (0,1) \\ 0 & \text{o/t} \end{cases}$

no special privilege to any $\theta \in \Theta$

$x = (0,1,1)$. Let's get best guess \downarrow consider as function of θ

$$\hat{\theta}_{MAP} = \underset{\theta \in (0,1)}{\operatorname{argmax}} \{P(\theta|x)\} = \underset{\theta \in (0,1)}{\operatorname{argmax}} \left\{ \frac{P(x|\theta)P(\theta)}{P(x)} \right\} = \underset{\theta \in (0,1)}{\operatorname{argmax}} \{P(x|\theta)\} = \underset{\theta \in (0,1)}{\operatorname{argmax}} \{\theta^2(1-\theta)\}$$

\uparrow
consider as function of θ

Same problem as solved before:

$$\frac{d}{d\theta} [\theta^2(1-\theta)] = \frac{d}{d\theta} (\theta^2 - \theta^3) = 2\theta - 3\theta^2 = 0 \Rightarrow \hat{\theta}_{MAP} = \frac{2}{3} = \hat{\theta}_{MLE}$$

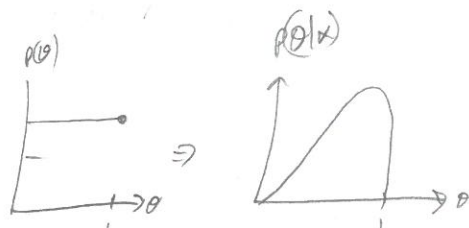
New question... What if we're interested in

$$P(\theta \in [0.6, 0.7] | X = (0, 1, 1))$$

How likely is the true θ between 0.6 and 0.7 after we see the data.

Frequentist answer: 0 or 1. Not possible!

Bayesian answer:
$$= \int_{0.6}^{0.7} P(\theta | X = (0, 1, 1)) d\theta$$



Let's try to solve for the posterior now.

$$P(\theta | x) = \frac{P(x|\theta) P(\theta)}{P(x)} = \frac{P(x|\theta)}{P(x)} = \frac{P(x|\theta)}{\int_0^1 P(x|\theta) P(\theta) d\theta} = \frac{P(x|\theta)}{\int_0^1 P(x|\theta) d\theta}$$

Bayes' theorem
(marginalizing over θ)

$$= \frac{\theta^2(1-\theta)}{\int_0^1 \theta^2(1-\theta) d\theta} = \frac{\theta^2(1-\theta)}{\int_0^1 (\theta^2 - \theta^3) d\theta} = \frac{\theta^2(1-\theta)}{\left[\frac{\theta^3}{3}\right]_0^1 - \left[\frac{\theta^4}{4}\right]_0^1} = \frac{\theta^2(1-\theta)}{\frac{1}{3} - \frac{1}{4}} = 12\theta^2(1-\theta)$$

$$\Rightarrow \int_{0.6}^{0.7} 12\theta^2(1-\theta) d\theta = 12 \left[\frac{\theta^3}{3} - \frac{\theta^4}{4} \right]_{0.6}^{0.7} = 0.1765 \quad \text{assuming the prior of indifference}$$

Let's solve for the posterior generally $P(\theta | x)$ where $X = (x_1, \dots, x_n)$ see 4 data pts.
and $\theta \sim U(0, 1)$

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)} = \frac{\prod_{i=1}^n P(x_i|\theta)}{\int P(x|\theta) P(\theta) d\theta} = \frac{\prod_{i=1}^n P(x_i|\theta)}{\int \prod_{i=1}^n P(x_i|\theta) d\theta}$$

(H) (0,1)

$$= \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}}{\int_0^1 \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} d\theta} = \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i}}{B(\sum x_i + 1, n - \sum x_i + 1)} = \text{Beta}(\sum x_i + 1, n - \sum x_i + 1)$$

Famous Integral, a "special domain"

this is a PDF of a
Binomial r.v

the Beta is
the posterior
of the iid
Bernoulli likelihood
and the uniform prior

$$B(\alpha, \beta) := \int_0^1 t^{\alpha-1} (1-t)^{\beta-1} dt$$

$$Y \sim \text{Beta}(\alpha, \beta) = \frac{1}{B(\alpha, \beta)} y^{\alpha-1} (1-y)^{\beta-1} \quad \text{Supp}[Y] = (0, 1)$$

Is this a valid PDF?

$$\int_0^1 \frac{1}{B(\alpha, \beta)} y^{\alpha-1} (1-y)^{\beta-1} dy = \frac{1}{B(\alpha, \beta)} \int_0^1 y^{\alpha-1} (1-y)^{\beta-1} dy = \frac{B(\alpha, \beta)}{B(\alpha, \beta)} = 1$$

Param space:

when is $\int_0^1 y^{\alpha-1} (1-y)^{\beta-1} dy$ finite?

If $\alpha = 0$ $y^{\alpha-1} = y^{-1}$ it's a hyperbola around $y=0 \Rightarrow$ diverge!

If $\beta = 0$ $(1-y)^{\beta-1} = (1-y)^{-1}$... $y=1 \Rightarrow$ diverge

the divergence is worse as α, β grow negative

$$\Rightarrow \boxed{\alpha > 0, \beta > 0}$$

$$E(Y) = \int_0^1 y \frac{1}{B(\alpha, \beta)} y^{\alpha-1} (1-y)^{\beta-1} dy = \frac{1}{B(\alpha, \beta)} \int_0^1 y^{\alpha} (1-y)^{\beta-1} dy$$

$$= \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)}$$

To simplify this, we need to introduce another famous integral:

$$\Gamma(\alpha) = \int_0^{\infty} t^{\alpha-1} e^{-t} dt \quad \text{the "gamma function"} \quad \forall \alpha > 0$$

we need two facts

① By integration by parts, you can show that

$$\Gamma(\alpha+1) = \alpha \Gamma(\alpha) \quad \text{the gamma function is the generalization of the factorial function, } \alpha!$$

$$\textcircled{2} B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)} \quad \text{which can be proven using ch. of var's or mult. var. calc.}$$

$$E(Y) = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \frac{\frac{\Gamma(\alpha+1) \Gamma(\beta)}{\Gamma(\alpha+1+\beta)}}{\frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}} = \frac{\frac{\alpha \Gamma(\alpha)}{(\alpha+\beta) \Gamma(\alpha+\beta)}}{\frac{\Gamma(\alpha)}{\Gamma(\alpha+\beta)}} = \boxed{\frac{\alpha}{\alpha+\beta}}$$

Var(Y) = ... similar (on HW)
+ tricks
applied

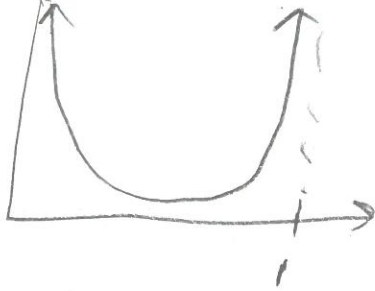
$$\text{Mode}(Y) = \underset{y \in (0,1)}{\operatorname{argmax}} \left\{ \frac{1}{B(\alpha, \beta)} y^{\alpha-1} (1-y)^{\beta-1} \right\} = \underset{y \in (0,1)}{\operatorname{argmax}} \left\{ y^{\alpha-1} (1-y)^{\beta-1} \right\} = \text{take log...}$$

$$= \underset{y \in (0,1)}{\operatorname{argmax}} \left\{ (\alpha-1) \ln(y) + (\beta-1) \ln(1-y) \right\} \Rightarrow \frac{\alpha-1}{y} - \frac{\beta-1}{1-y} = 0$$

$$\Rightarrow \frac{1-y}{y} = \frac{\beta-1}{\alpha-1} \Rightarrow \frac{1}{y} - 1 = \frac{\beta-1}{\alpha-1} \Rightarrow \frac{1}{y} = \frac{\beta-1+\alpha-1}{\alpha-1} \Rightarrow \boxed{\text{Mode}(Y) = \frac{\alpha-1}{\alpha+\beta-2}}$$

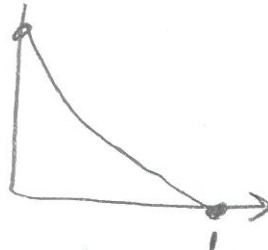
2nd deriv. check on HW

Shapes of the beta r.v. density:

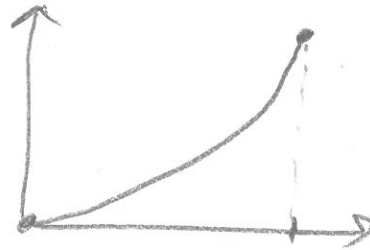


$$\alpha = \beta = 0.5$$

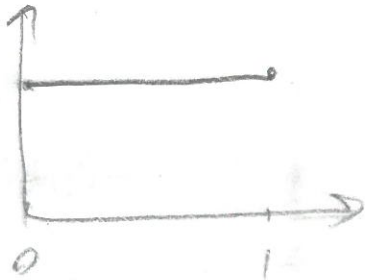
"arcsin distribution"



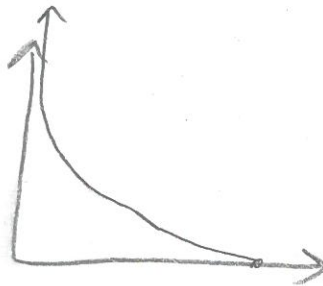
$$\alpha = 1, \beta = 3$$



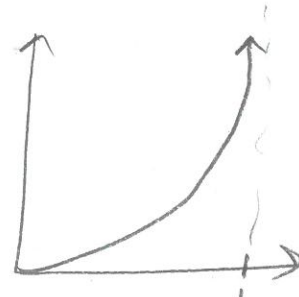
$$\alpha = 5, \beta = 1$$



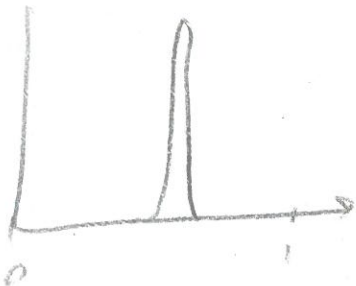
$$\alpha = \beta = 1$$



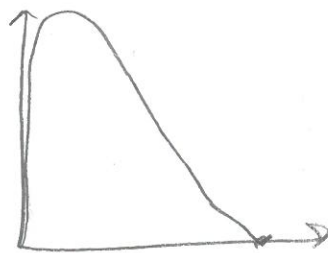
$$\alpha = 0.99, \beta = 3$$



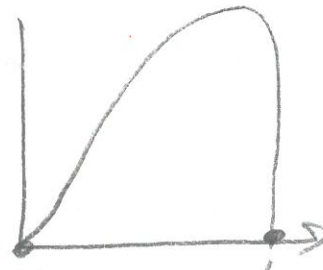
$$\alpha = 5, \beta = 0.99$$



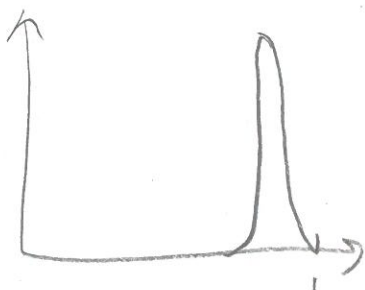
$$\alpha = \beta = 100$$



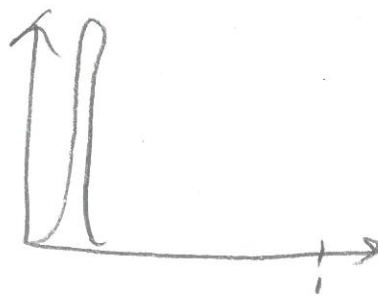
$$\alpha = 1.01, \beta = 3$$



$$\alpha = 5, \beta = 1.01$$



$$\alpha = 100, \beta = 10$$



$$\alpha = 10, \beta = 100$$