

- Lecture 6

Posterior predictive distribution - finding  $X_4$  from  $X_1, X_2, X_3$

$\theta \sim P(\theta)$  - prior

$X_1, X_2, X_3$  - data observed

we want to find:

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

Note:  $P(Y) = \sum_{x \in \text{supp}(X)} P(X, Y)$  so:

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

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

Note:  $P(X, \theta) = P(X | \theta) P(\theta)$  so:

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

\*Note:  $P(X_4 | \theta, X_1, X_2, X_3) = P(X_4 | \theta)$

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

so:

- $$P(X_4 | X_1, X_2, X_3) = \sum_{\theta \in \Theta_0} P(X_4 | \theta) \underbrace{P(\theta | X_1, X_2, X_3)}_{\text{posterior prior}}$$

Note:

$$\neq P(X_4 | \hat{\theta}_{MLE})$$



• Having  $\theta$  range over all  $\Theta = (0, 1)$

example  $\Theta_0 = \{0.25, 0.75\}$   $X_1, X_2, X_3$  observed  
 $X = \langle 0, 1, 1 \rangle$

$$P(\theta) = \bigcup_{\text{uniform}} P(\Theta_0) = \begin{cases} 0.25 \text{ w.p. } \frac{1}{2} \\ 0.75 \text{ w.p. } \frac{1}{2} \end{cases}$$

• 1st Goal: Point estimation (improved)

$$\hat{\theta}_{\text{MAP}} := \arg \max_{\theta \in \Theta_0} \{P(\theta|X)\} = \arg \max_{\theta \in \Theta_0} \left\{ \frac{P(X|\theta)P(\theta)}{P(X)} \right\} = \arg \max_{\theta \in \Theta_0} \{P(X|\theta)\} = \arg \max_{\theta \in \Theta_0} \{f(\theta; X)\} = \hat{\theta}_{\text{MLE}}$$

constant if prior is uniform

constant in  $\theta$

\* This ignores possible values for  $\theta$  in the space \*

• for ex: it was shown  $\hat{\theta}_{\text{MAP}} = 0.75 \neq \hat{\theta}_{\text{MLE}} = 0.66$ . B/C  $\Theta_0 \neq \Theta$

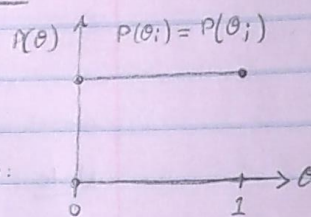
We want:

\* Parameter space of the likelihood model = support of the prior

$$[\Theta = (0, 1) = \text{supp}[\theta] = (0, 1)]^*$$

$$P(\theta) = \begin{cases} 1 & \text{if } \theta \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

\* uniform, so:



• for ex: we want given  $X_1, X_2, X_3$ :  $P(\theta): \theta \sim U(0, 1)$  \*\* standard uniform r.v

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta \in \Theta_0} \left\{ \frac{P(X|\theta)P(\theta)}{P(X)} \right\} = \arg \max_{\theta \in \Theta_0} \{P(X|\theta)\} = \arg \max_{\theta \in (0, 1)} \left\{ \theta^2(1-\theta) \right\}$$

constant = 1

constant for  $\theta \in \Theta_0$

$\hat{x} = 1$   $\hat{x} = 0$

\* Note to find max set  $\frac{d}{d\theta}[P(X|\theta)] = 0$

$$\text{so: } \frac{d}{d\theta}[\theta^2(1-\theta)] = \frac{d}{d\theta}[\theta^2 - \theta^3] = 2\theta - 3\theta^2 = 0 \Rightarrow 2\theta = 3\theta^2 \Rightarrow 2 = 3\theta$$

$$\text{Now } \hat{\theta}_{\text{MAP}} = \frac{2}{3} = \hat{\theta}_{\text{MLE}}$$

$$\theta = \frac{2}{3}$$

\* frequentism this isn't possible, only 1 (in range) or 0 (not in range)

• 2nd Goal: Probability of falling in some range.  $P(\theta \in [0.6, 0.7] | X = \langle 0, 1, 1 \rangle)$

$$\begin{aligned} \int_{0.6}^{0.7} P(\theta|X) d\theta &= \int_{0.6}^{0.7} \frac{P(X|\theta)P(\theta)}{P(X)} d\theta = \int_{0.6}^{0.7} \frac{\theta^2(1-\theta) \cdot 1}{\int_{\Theta} P(X|\theta)P(\theta) d\theta} d\theta \\ &= \int_{0.6}^{0.7} \frac{\theta^2 - \theta^3}{\frac{1}{12}} d\theta = 12 \left[ \frac{\theta^3}{3} - \frac{\theta^4}{4} \right]_{0.6}^{0.7} \approx 17.65\% \end{aligned}$$

chance  $\theta$  is between  $[0.6, 0.7]$



Generalization: Arbitrary data set

$n$  realizations of the iid Bernoulli model:  $X_1, X_2, \dots, X_n$

$$\text{for single } X \quad P(\theta|X) = \frac{P(X|\theta) P(\theta)}{P(X)} = \frac{P(X|\theta) \cdot 1}{\int_{\Theta} P(X|\theta) P(\theta) d\theta} = \frac{P(X|\theta)}{\int_{\Theta} P(X|\theta) d\theta}$$

$$\text{for } X_i \quad \{i=1,2,\dots,n\} \quad P(\theta|X_1, X_2, \dots, X_n) = \frac{\prod_{i=1}^n P(X_i|\theta)}{\int_{\Theta} \prod_{i=1}^n P(X_i|\theta) d\theta} = \frac{\theta^{\sum X_i} (1-\theta)^{n-\sum X_i}}{\int_0^1 \theta^{\sum X_i} (1-\theta)^{n-\sum X_i} d\theta}$$

• Note

where  $X_i$  is either 1 or 0, so  $\sum X_i$  is all 'successes' and  $n - \sum X_i$  is 'failures'

$$= \frac{\theta^{\sum X_i} (1-\theta)^{n-\sum X_i}}{\int_0^1 \theta^{\sum X_i + 1} (1-\theta)^{n-\sum X_i + 1} d\theta} \quad \left. \vphantom{\int_0^1} \right\} \text{Beta function}$$

• Note

Beta functions are:

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

$$** P(\theta|X) = \frac{\theta^{\sum X_i + 1} (1-\theta)^{n-\sum X_i + 1}}{B(\sum X_i + 1, n - \sum X_i + 1)} := \left[ \text{Beta}(\sum X_i + 1, n - \sum X_i + 1) \right]$$

PDF for Beta distribution

• Properties of the Beta distribution

$$\text{let } Y \sim \text{Beta}(\alpha, \beta) := \frac{1}{B(\alpha, \beta)} \cdot y^{\alpha-1} (1-y)^{\beta-1} = f(y)$$

•  $\text{supp}[Y] = (0, 1)$

• parametric space:

$$\alpha > 0 \\ \beta > 0$$

is  $f(y)$  a PDF? does  $\int_{\text{supp}[Y]} f(y) dy = 1$ ?

$$\int_{\text{supp}[Y]} f(y) dy = \int_0^1 \frac{1}{B(\alpha, \beta)} \cdot y^{\alpha-1} (1-y)^{\beta-1} dy = \frac{1}{B(\alpha, \beta)} \int_0^1 y^{\alpha-1} (1-y)^{\beta-1} dy = \frac{B(\alpha, \beta)}{B(\alpha, \beta)} = 1$$

so  $f(y)$  is a PDF

• Note: Gamma function

$$\Gamma(\alpha) := \int_0^{\infty} t^{\alpha-1} e^{-t} dt$$

for  $\forall \alpha > 0$

properties

$$\textcircled{a} \quad \Gamma(\alpha+1) = \alpha \Gamma(\alpha)$$

(integration by parts)

$$\textcircled{b} \quad B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

(mult. var. calc.)



- Note  $\alpha = \text{success} + 1$   
 $\beta = n - (\text{success} + 1)$

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

use properties  
of gamma function

$$\propto \Gamma(\alpha)$$

• Expectation  $E[Y] = \int_{\text{supp}[Y]} y f(y) dy = \int_0^1 y \cdot \frac{y^{\alpha-1} (1-y)^{\beta-1}}{B(\alpha, \beta)} dy = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \frac{\frac{\Gamma(\alpha+1) \Gamma(\beta)}{\Gamma(\alpha+\beta)}}{\frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)}} = \frac{\Gamma(\alpha)}{\Gamma(\alpha+\beta)} \cdot \frac{\Gamma(\alpha+1) \Gamma(\beta)}{\Gamma(\alpha) \Gamma(\beta)} = \frac{\alpha}{\alpha+\beta}$

$$E[Y] = \frac{\alpha}{\alpha+\beta}$$

• Mode  $\text{Mode}[Y] = \underset{y \in \text{supp}[Y]}{\text{argmax}} \{f(y)\} = \underset{y \in (0,1)}{\text{argmax}} \left\{ \frac{1}{B(\alpha, \beta)} \cdot y^{\alpha-1} (1-y)^{\beta-1} \right\} = \underset{y \in (0,1)}{\text{argmax}} \left\{ y^{\alpha-1} (1-y)^{\beta-1} \right\}$

(greatest Y value) constant

$$= \underset{y \in (0,1)}{\text{argmax}} \{(\alpha-1) \ln(y) + (\beta-1) \ln(1-y)\}$$

• Note: find argmax by  $\frac{d}{dy} [f(y)] = 0$

$$\Rightarrow \frac{d}{dy} [\text{" "}] = \frac{\alpha-1}{y} - \frac{\beta-1}{1-y} = 0$$

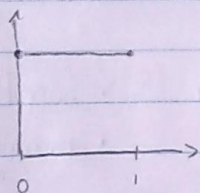
$$\Rightarrow \frac{\alpha-1}{y} = \frac{\beta-1}{1-y}$$

$$\Rightarrow y = \frac{\alpha-1}{\alpha+\beta-2} \quad \text{when } \alpha > 1, \beta > 1 = \text{Mode}[Y]$$

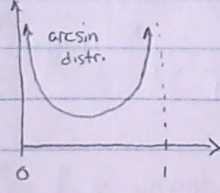
• Beta r.v. density graphed:

$f(y)$   
 $y$

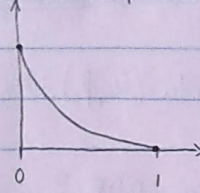
$\alpha = \beta = 1 : U(0,1)$



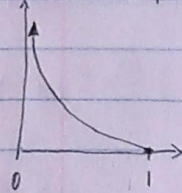
$\alpha = \beta = 0.5$



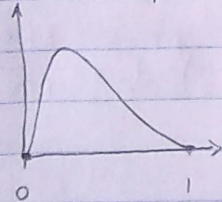
$\alpha = 1, \beta = 3$



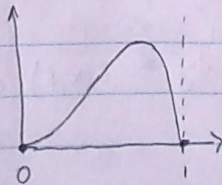
$\alpha = 0.99, \beta = 3$



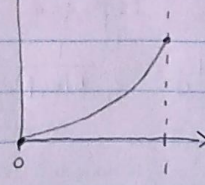
$\alpha = 1.01, \beta = 3$



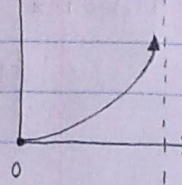
$\alpha = 3, \beta = 1.01$



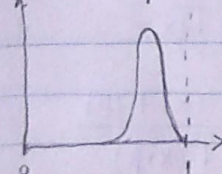
$\alpha = 3, \beta = 1$



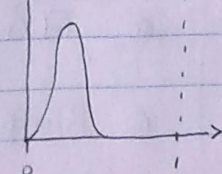
$\alpha = 3, \beta = 0.99$



$\alpha = 100, \beta = 10$



$\alpha = 10, \beta = 100$



$\alpha = \beta = 100$

