

## 1. Identifying and matching kernels.

Part I: Refresher: kernel = taking out constants that don't involve variable

a)  $f(\theta | y) = y e^{-\theta y} : \theta > 0$

$$\Rightarrow \text{kernel: } e^{-\theta y}$$

b)  $\frac{\theta^y}{(y-1)!} = \text{scalar} \rightarrow \text{kernel: } \theta^{y-1} e^{-2\theta}$

c) scalar = 3, kernel =  $\theta^2$

Part II: Refresher: Beta(a, b) =  $\theta^{a-1} (1-\theta)^{b-1}$

a)  $f(\theta | y) \propto \text{Beta}(3, 12)$

b)  $f(\theta | y) \propto \theta^{12-1} (1-\theta)^{3-1} = \text{Beta}(12, 3)$

c)  $f(\theta | y) \propto \theta^{1-1} (1-\theta)^{1-1} = \text{Beta}(1, 1)$

## 2. Bernoulli:

a) The joint density of  $y_1, y_2, \dots, y_n$  is the product of all the  $y_i$ 's since they are iid.

$$\Rightarrow p(y_1, y_2, \dots, y_n | \theta) = p(y_1 | \theta) \cdot p(y_2 | \theta) \cdots p(y_n | \theta)$$

$$\text{for } p(y_i | \theta) = \theta^{y_i} (1-\theta)^{1-y_i}$$

$$\Rightarrow p(y_1, y_2, \dots, y_n | \theta) = \prod_{i=1}^n (\theta^{y_i} \cdot (1-\theta)^{1-y_i})$$

b) Prior:  $\text{Beta}(a, b) \Rightarrow p(\theta) = \frac{\theta^{a-1} (1-\theta)^{b-1}}{B(a, b)}$   $\rightarrow$  beta function

$$p(y | \theta) = \theta^y (1-\theta)^{1-y}$$

Using Bayes' Rule:

$$p(\theta | y_1, y_2, \dots, y_n) = \frac{p(y_1, y_2, \dots, y_n | \theta) \cdot p(\theta)}{p(y_1, y_2, \dots, y_n)}$$

$$p(y_1, y_2, y_3, \dots, y_n) = \underbrace{\int_0^1}_{\text{everything else}} p(y_1, y_2, \dots, y_n | \theta) \cdot p(\theta) d\theta$$

$$= \int_0^1 \prod_{i=1}^n (\theta^{y_i} \cdot (1-\theta)^{1-y_i}) \cdot \frac{\theta^{a-1} (1-\theta)^{b-1}}{B(a, b)} d\theta$$

$$= \frac{1}{B(a, b)} \underbrace{\int_0^1}_{\text{everything else}} \dots$$

$$\Rightarrow p(\theta | y_1, y_2, \dots, y_n) = \frac{\prod_{i=1}^n (\theta^{y_i} \cdot (1-\theta)^{1-y_i}) \cdot \theta^{a-1} (1-\theta)^{b-1}}{B(a, b) \cdot \int_0^1 \prod_{i=1}^n (\theta^{y_i} \cdot (1-\theta)^{1-y_i}) d\theta}$$

Let  $k = \sum_{i=1}^n y_i$ . Also consider the fact that the sum of independent

Bernoulli variables follows a binomial distribution:

$p(a) = \binom{n}{k} p^k \cdot (1-p)^{n-k}$  (for  $n, p, k$  unrelated to any values in our problem. This is just the formula). Consider again:

$$p(\theta | y_1, y_2, \dots, y_n) = \frac{\text{numerator } \textcircled{*}}{\text{denominator}}$$

$$= \frac{\text{numerator } \textcircled{*}}{\int_0^1 \binom{n}{k} \theta^k \cdot (1-\theta)^{n-k} d\theta}$$

$$= \text{numerator } \textcircled{*}$$

$$= \theta^{(\sum y_i + a) - 1} \cdot (1-\theta)^{(n - \sum y_i + b) - 1}$$

$$= \text{Beta}(\sum y_i + a, n - \sum y_i + b)$$

this is just a scalar so it won't affect the posterior distribution's form

3.

$$\text{Bayes' Rule: } p(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)}$$

fixed  
calculations  
errors  
thanks  
to peer  
reviews

$$\Rightarrow p(y|\theta) \cdot Y|\theta \sim \text{Bin}(6, \theta)$$

$$\Rightarrow p(y|\theta) = \binom{6}{y} \theta^y (1-\theta)^{6-y}, \text{ Kasparov won 1 game}$$

$$\Rightarrow p(\theta|y) = \binom{6}{y} \theta^y (1-\theta)^{6-y} / 6\theta(1-\theta)^5$$

↳ Calculating posteriors:

$$\begin{aligned} - t = 0.2 &\Rightarrow p(\theta = 0.2 | y=1) \propto p(y|0.2) \cdot p(0.2) \\ &= 6 \times 0.2 \times (1-0.2)^5 \cdot 0.1 \approx 0.039 \end{aligned}$$

$$\begin{aligned} - \theta = 0.5 &\propto 6 \times 0.5 \times (1-0.5)^5 \times 0.25 \\ &\approx 0.023 \end{aligned}$$

$$\begin{aligned} - \theta = 0.8 &\propto 6 \times 0.8 \times (1-0.8)^5 \times 0.65 \\ &\approx 0.000998 \end{aligned}$$

$$\Rightarrow p(y) = \sum_{\theta} p(y|\theta) \cdot p(\theta)$$

$$\begin{aligned} &= p(y|0.2) \cdot p(0.2) + p(y|0.5) \cdot p(0.5) + p(y|0.8) \cdot p(0.8) \\ &\approx 0.088 \end{aligned}$$

$$\Rightarrow p(\theta = 0.2 | y=1) \approx \frac{0.039}{0.088} \approx 0.44$$

$$p(\theta = 0.5 | y=1) \approx \frac{0.023}{0.088} \approx 0.26$$

$$p(\theta = 0.8 | y=1) \approx \frac{0.000998}{0.088} \approx 0.011$$

$$4. E[L(\theta, \hat{\theta}) | y] = \int L(\theta, \hat{\theta}) \cdot p(\theta | y) d\theta \quad \begin{matrix} (\text{because the expectation}) \\ (\text{is taken with respect}) \\ (\text{to } \theta | y) \end{matrix}$$

$$\begin{aligned} \text{refresher: } E[\theta | y] &= \int \theta \cdot p(\theta | y) d\theta \\ &= \int (\theta^2 - 2\theta\hat{\theta} + \hat{\theta}^2) \cdot p(\theta | y) d\theta \\ &= \int \theta^2 \cdot p(\theta | y) d\theta - 2\theta\hat{\theta} \cdot p(\theta | y) d\theta + \hat{\theta}^2 \cdot p(\theta | y) d\theta \\ &= \left( \int \theta^2 \cdot p(\theta | y) d\theta \right) - \left( \int 2\theta\hat{\theta} \cdot p(\theta | y) d\theta \right) + \left( \int \hat{\theta}^2 \cdot p(\theta | y) d\theta \right) \\ &= E[\theta^2 | y] - E[2\theta\hat{\theta} | y] + E[\hat{\theta}^2 | y] \end{aligned}$$

$$(\text{but } \hat{\theta} \text{ is a constant}) = E[\theta^2 | y] - 2\hat{\theta} \cdot E[\theta | y] + \hat{\theta}^2$$

↳ this is minimized when

derivative with respect to  $\hat{\theta}$ , which is not a constant here since we need to find a value for it

$$\frac{d}{d\hat{\theta}} \left( E[\hat{\theta}^2 | y] - 2\hat{\theta} E[\hat{\theta} | y] + \hat{\theta}^2 \right) = 0$$

$$\Rightarrow -2E[\hat{\theta} | y] + 2\hat{\theta} = 0 \Rightarrow \hat{\theta} = E[\hat{\theta} | y]$$

5. a) General  $a$  and  $b$  choosing guidelines:

$$E(\theta) = \frac{a}{a+b}. \text{ Our mean: } 15\%$$

Higher  $a$  and  $b$  means lower variability. Our range:  $10\% - 25\%$

Based on these, we can use the model  $\hat{\theta} \sim \text{Beta}(2, 11)$

$$b) a_{\text{post}} = a_{\text{prior}} + \text{# successes} = 2 + 30 = 32$$

$$b_{\text{post}} = b_{\text{prior}} + \text{# fails} = 11 + 60 = 71$$

$$\Rightarrow \hat{\theta} | y \sim \text{Beta}(32, 71) \Rightarrow y=30, n=90$$

$$c) E[\hat{\theta} | y] = \frac{a+y}{a+b+n} = \frac{32+30}{32+11+90} \approx 0.31$$

$$\text{Posterior mode} = \frac{a_{\text{post}} - 1}{a_{\text{post}} + b_{\text{post}} - 2} = \frac{31}{101} \approx 0.307$$

$$\text{Posterior SD} = \sqrt{\frac{a_{\text{post}} \times b_{\text{post}}}{(a_{\text{post}} + b_{\text{post}})^2 (a_{\text{post}} + b_{\text{post}} + 1)}} = \sqrt{\frac{32 \times 71}{(32+71)^2 (32+71+1)}} \approx 0.045$$

- d) The model reflects the data more:

- Posterior mean  $\approx 31\%$ , which is close to the proportion of success in the data.

- The prior model sample is much smaller than the data sample, so the posterior mean is calculated with more weight on the data success rate.

$$b. p(\lambda | y_1, y_2, \dots, y_n) \propto p(y_1, y_2, \dots, y_n | \lambda) \cdot p(\lambda)$$

$$\text{Refresher: } p(y_i | \lambda) = \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} \text{ for Poisson}$$

$(y_1, y_2, \dots, y_n \text{ iid})$

$$\begin{aligned} \Rightarrow p(y_1, y_2, \dots, y_n | \lambda) &= p(y_1 | \lambda) \cdot p(y_2 | \lambda) \cdots p(y_n | \lambda) \\ &= \prod_{i=1}^n \left( \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} \right) = \frac{e^{-n\lambda} \lambda^{\sum y_i}}{\prod_{i=1}^n y_i!} \end{aligned}$$

$$\Rightarrow p(\lambda | y_1, y_2, \dots, y_n) \propto e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n y_i}}{\prod_{i=1}^n y_i!} \cdot \frac{b^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-b\lambda}$$

$$\propto \frac{b^\alpha}{\Gamma(\alpha)} e^{-n\lambda} \lambda^{\sum_{i=1}^n y_i + \alpha - 1} \cdot e^{-(b+1)\lambda}$$

Refresher:  $f(x | k, \theta) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-\frac{x}{\theta}}$  (gamma PDF)

$$\Rightarrow p(\lambda | y_1, y_2, \dots, y_n) = \text{Gamma}(\sum_{i=1}^n y_i + \alpha, b+1)$$

because!  $\lambda$ : variable  
 $\Rightarrow \lambda^{\sum_{i=1}^n y_i + \alpha - 1}$  gives  $\sum y_i + \alpha$  for the shape

$$\Rightarrow e^{-(b+1)\lambda} = e^{-\frac{\lambda}{\theta}} \Rightarrow -(b+1)\lambda = -\theta \frac{\alpha}{\theta} \lambda$$

$$\Rightarrow \theta = b+1$$

7. In our context,  $\lambda$  is the random variable

- Parameter  $\alpha$ : decides the shape. The bigger  $\alpha$  is, the smaller the interval on the curve (Beta is also like this). This also implies that the prior is strong and might weigh more in calculating the posterior. Bigger  $\alpha$  also indicates less variability. The reverse is true for a small parameter  $\alpha$ . In a Poisson context,  $\alpha$  represents the number of occurrences in an interval.

- Parameter  $b$ : this is the exponent of  $e$ . The larger  $b$  is, the more concentrated the distribution is. This means lower variability in the random variable  $\lambda$ .