

Lec 4 Math 341 2/7/19

Let  $\theta$  be a r.v. then this makes sense:

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

likelihood

prior: your thoughts about  $\theta$  before you see any data

posterior:  
thoughts about  
 $\theta$  after you  
see data

prior predictive distr.

$$p(x) = \sum_{\theta \in \Theta} p(x|\theta) p(\theta)$$

$$p(x) = \int p(x|\theta) p(\theta) d\theta$$

a r.v. with a value  
continuous or  
a count  
in  $\Theta$

$$p(x|\theta) = p(x; \theta)$$

$$= L(\theta; x)$$

example

$F = \text{ich familli}$ ,  $x = \langle 0, 1, 1 \rangle$   $p(x|\theta) = \theta^2(1-\theta)$

let  $\Theta_0 = \{0.5, 0.75\}$  i.e. not the full parameter space.

$p(\theta = 0.75|x) > p(\theta = 0.5|x)$  ? Let's see how it works...

$$p(\theta = 0.75|x) = \frac{p(x|\theta = 0.75) p(\theta = 0.75)}{p(x|\theta = 0.75) p(\theta = 0.75) + p(x|\theta = 0.5) p(\theta = 0.5)}$$

$$p(x|\theta = 0.75) = 0.75^2 \cdot 0.25 = .141$$

$$p(x|\theta = 0.5) = 0.5^2 \cdot 0.5 = 0.125$$

we need  $p(\theta = 0.75)$  &  $p(\theta = 0.5)$ . That is  $p(\theta)$ ! It is subjective! What do you think??

# Principle of Indifference

all  $\theta \in \Theta$  equally likely a priori i.e.  $P(\theta) = \frac{1}{|\Theta|}, \forall \theta$

Here...

if  $\theta$  is discrete

$$P(\theta) = \begin{cases} 0.5 & \text{if } \theta = 0.75 \\ 0.5 & \text{if } \theta = 0.5 \end{cases}$$

$$P(\theta = 0.75 | x) = \frac{(.141)(0.5)}{(.141)(0.5) + (.125)(0.5)} = 0.53$$

$$P(\theta = 0.5 | x) = \frac{.125}{.141 + .125} = 0.47$$

After the data, the more likely param value is 0.75.

$$P(\theta = 0.75) = 0.5 \xrightarrow{x} P(\theta = 0.75 | x) = 0.53$$

## Bayesian Conditionalization

Let's look at the full space: all  $x$  and all  $\theta$

$$x \in \mathcal{X} = \{0,1\} \times \{0,1\} \times \{0,1\}$$

$\mathcal{X}$

								$\langle 0,0,0 \rangle$	
$\langle 1,1,1 \rangle$				$\langle 1,1,0 \rangle$	$\langle 1,0,1 \rangle$	$\langle 0,1,1 \rangle$	$\langle 0,0,1 \rangle$	$\langle 0,1,0 \rangle$	0.75
								$\Theta$	
$\langle 1,1,1 \rangle$	$\langle 1,1,0 \rangle$	$\langle 1,0,1 \rangle$	$\langle 0,1,1 \rangle$	$\langle 1,0,0 \rangle$	$\langle 0,1,0 \rangle$	$\langle 0,0,1 \rangle$	$\langle 0,0,0 \rangle$	0.5	

$$P(X=(1,1,1) | \theta=.75) = .922$$

$$P(X=(1,1,0) | \theta=.75) = .141$$

$$P(X=(1,0,0) | \theta=.75) = .047$$

$$P(X=(0,0,0) | \theta=.75) = .016$$

Is  $\theta$  indep of  $X$ ? NO

$X$  tells you something about  $\theta$ ...

that's the whole pt. of this course!

Cool questions:

$$P(X=(1,0,0), \theta=0.5) = .125 \cdot .5 = .0625$$

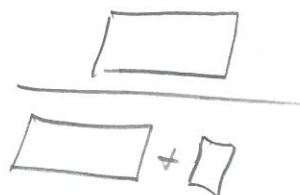


$$P(X=(1,0,0)) = P(X=(1,0,0), \theta=.75) + P(X=(1,0,0), \theta=.5)$$

$$= .047 \cdot .5 + .125 \cdot .5 = .086$$



$$P(\theta=0.5 | X=(1,0,0)) = \frac{.0625}{.086} = 0.73$$



$$\sum_{\theta \in \Theta} P(\theta) = 1$$

$$\sum_{\theta \in \Theta} P(\theta | x) = 1$$

$$\sum_{\theta \in \Theta} P(x | \theta) = ?$$

not guaranteed  
to be  
satisfying

See above  $x=(1,1,1)$

$$\Theta_0 = \{0.1, 0.25, .5, 0.75, 0.9\}$$

$$P(\theta) = \{0.2 \text{ for } \forall \theta \in \Theta\}$$

always always a constant is  $\theta$   
under principle of indifference

$$P(\theta | x) = \frac{1}{P(x)} P(x | \theta) P(\theta) \propto P(x | \theta) P(\theta) \propto P(x | \theta)$$

view this as a  
function of  $\theta$  Same for all  $\theta$

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

scale based on prior  
scale vs rest of data  
relief

$$P(X|0=0.1) = .1^2 \cdot .9 = .009$$

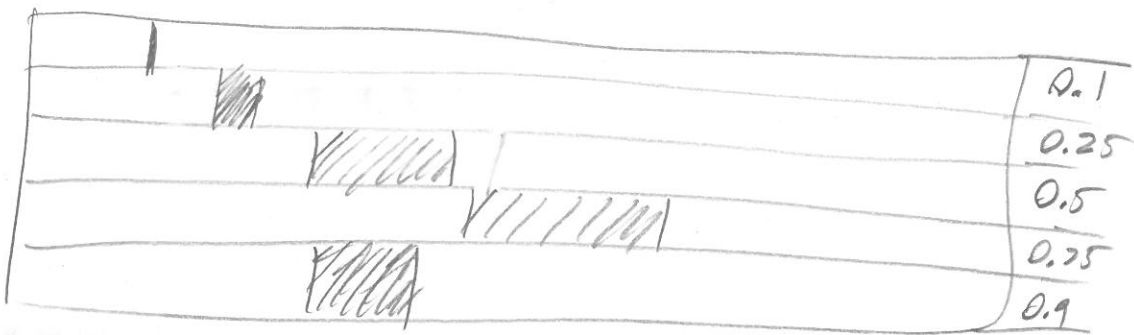
$$P(X|0=0.25) = .25^2 \cdot .75 = .047$$

$$P(X|0=0.5) = .125$$

$$P(X|0=.75) = .141$$

$$P(X|0=.9) = .9^2 \cdot .1 = .081$$

$X$



If we want the most likely value of  $\theta$  given the data,

what is it?  $\theta = 0.75$

always true since  $P(X)$  constant  $\theta$

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

Maximum a posteriori

(Maximum mode)

a Bayesian point estimate!

True if principled inference

$$= \hat{\theta}_{MLE}$$

but  $\hat{\theta}_{MLE} = 0.66 \neq 0.75$  What happened?

not true  $\hat{\theta}_{MLE}$  requires search over full  $\Theta$ , we're not there yet.

Let's return to  $\Theta_0 = \{0.5, 0.75\}$ ,  $X = (0, 1, 1)$

And let's look at our  $P(\Theta = 0.75)$  iteratively...

$$P(\Theta = 0.75 | X_1 = 0) = \frac{P(X_1 = 0 | \Theta = 0.75) \cdot P(\Theta = 0.75)}{P(X_1 = 0 | \Theta = 0.75) P(\Theta = 0.75) + P(X_1 = 0 | \Theta = 0.5) P(\Theta = 0.5)}$$

Use this  
as prior

$$= \frac{\frac{0.25}{0.75 + 0.5} = \frac{1}{3}}{\frac{1}{3}}$$

$$P(\Theta = 0.75 | X_2 = 1) = \frac{P(X_2 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_1 = 0)}{P(X_2 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_1 = 0) + P(X_2 = 1 | \Theta = 0.5) P(\Theta = 0.5 | X_1 = 0)}$$

$$P(X_2 | X_1) = \sum_{\Theta \in \Theta_0} P(X_2 | \Theta, X_1) P(\Theta | X_1)$$

$$\frac{0.75}{0.75} \quad \frac{1}{3} \quad \frac{0.5}{0.5} \quad \frac{2}{3}$$

$$= 0.429 \quad , 429$$

$$P(\Theta = 0.75 | X_3 = 1) = \frac{P(X_3 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_1 = 0, X_2 = 1)}{P(X_3 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_1 = 0, X_2 = 1) + P(X_3 = 1 | \Theta = 0.5) P(\Theta = 0.5 | X_1 = 0, X_2 = 1)}$$

$$\frac{0.75}{0.75} \quad , 429 \quad \frac{0.5}{0.5} \quad , 571$$

$$= 0.53$$

$$P(\Theta = 0.75) = 0.5 \xrightarrow{X_1} 0.33 \xrightarrow{X_2} 0.429 \xrightarrow{X_3} 0.53$$

constant refinement

WTS  $P(\theta | x_1, \dots, x_n) = \frac{P(x_n | \theta) P(\theta | x_1, \dots, x_{n-1})}{\sum_{\theta \in \Theta} P(x_n | \theta) P(\theta | x_1, \dots, x_{n-1})}$  H<sub>3</sub>

Start with  $P(\theta | x_1, \dots, x_n) = \frac{P(x_1, \dots, x_n, \theta | \theta) P(\theta)}{P(x_1, \dots, x_{n-1}, x_n)}$  ind. r.v.'s

$= \frac{P(x_1 | \theta) \cdot \dots \cdot P(x_{n-1} | \theta) P(x_n | \theta) P(\theta)}{P(x_n | x_1, \dots, x_{n-1}) P(x_1, \dots, x_n)}$  def of cond prob

$= \frac{P(x_n | \theta) P(x_1, \dots, x_{n-1} | \theta) P(\theta)}{P(x_n | x_1, \dots, x_{n-1}) P(x_1, \dots, x_n)} P(\theta | x_1, \dots, x_{n-1})$

$$P(x_n | x_1, \dots, x_{n-1}) = \sum_{\theta \in \Theta} P(x_n, \theta | x_1, \dots, x_{n-1})$$

Why? I can introduce another r.v. and margin over it.

$$P(x) = \sum_{y \in \mathcal{Y}} P(x, y)$$

Then I can condition on a third r.v. Z

$$P(x | z) = \sum_{y \in \mathcal{Y}} P(x, y | z)$$

$$= \sum_{\theta \in \Theta} P(x_n | \theta, x_1, \dots, x_{n-1}) P(\theta | x_1, \dots, x_{n-1})$$

7

$$\begin{aligned}
 \text{Note } P(X_n | \theta, X_1, \dots, X_{n-1}) &= \frac{P(X_1, \dots, X_{n-1}, X_n, \theta)}{P(X_1, \dots, X_{n-1}, \theta)} \\
 &= \frac{P(X_1, \dots, X_{n-1}, X_n | \theta) \cancel{P(\theta)}}{P(X_1, \dots, X_{n-1} | \theta) \cancel{P(\theta)}} \\
 &= \frac{\cancel{P(X_1 | \theta)} \dots \cancel{P(X_{n-1} | \theta)} P(X_n | \theta)}{\cancel{P(X_1 | \theta)} \dots \cancel{P(X_{n-1} | \theta)}} = P(X_n | \theta)
 \end{aligned}$$

And can you see this?

If  $\theta$  is known, then  $X_1, \dots, X_{n-1}$  do not give any more information to the distr. of  $X_n$ .  $\theta$  completely governs the distr.

---