

lec 4 2/17/17 Prob 341

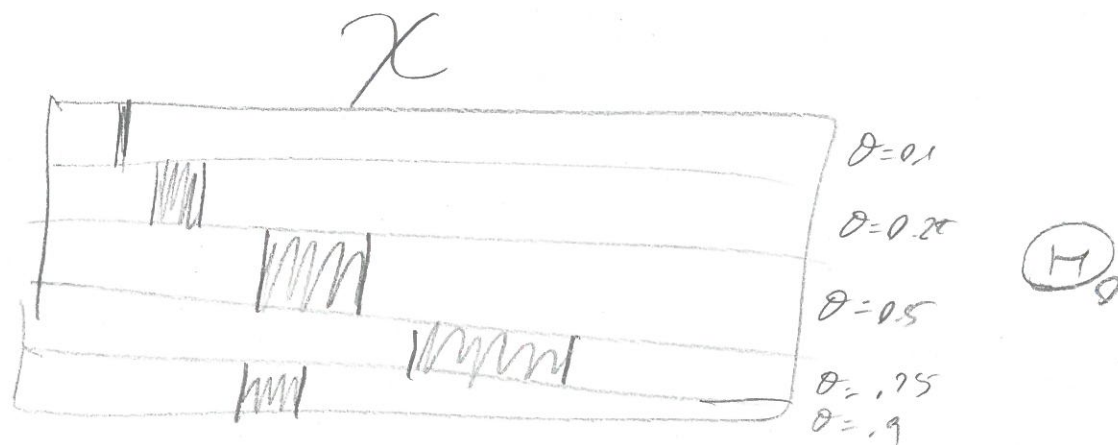
$F = \text{Bernoulli}$

$X = (0, 1, 1)$

$$\Theta_0 = \{0.1, 0.25, 0.5, 0.75, 0.9\}$$

$\theta \sim U(\Theta_0)$ i.e. discrete uniform on the elements of Θ_0 i.e. $P(\theta) = 0.2$

We want to find $P(\theta|x)$. Before we do that... draw picture



$P(x|\theta)$ represents the prop. of area in any slice.

$$P(x|\theta=0.1) = 0.009$$

$$P(x|\theta=0.25) = 0.047$$

$$P(x|\theta=0.5) = 0.125$$

$$P(x|\theta=0.75) = 0.191$$

$$P(x|\theta=0.9) = 0.061$$

$P(\theta|x)$ is the slices ^{each of} div by the total area of all slices

What is best model already? Biggest slice of the slices $\theta = 0.75$

Already you see its $P(\theta|x)$ is going to be the largest

Bigger slice of the slices is a form of pt. estimation (over guess of θ)! [2]

$$\hat{\theta}_{\text{MAP}} = \underset{\theta \in \mathcal{H}_0}{\text{argmax}} \{ P(\theta|x) \} = \underset{\theta \in \mathcal{H}_0}{\text{argmax}} \left\{ \frac{P(x|\theta)P(\theta)}{P(x)} \right\}$$

Max. a posteriori
Bayesian estimate
(AKA "posterior mode")

$$= \underset{\theta \in \mathcal{H}_0}{\text{argmax}} \{ P(x|\theta) P(\theta) \}$$

since $P(x)$ is
a normalizing
constant
 $\neq f(\theta)$

$$= \underset{\theta \in \mathcal{H}_0}{\text{argmax}} \{ P(x|\theta) \}$$

since $P(\theta)$ is
the same for all \mathcal{H}_0
 $\neq f(\theta)$

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

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

scale factor
on prior
belief
(height)

normalize so
all $P(\theta|x)$'s add
up to 1
(relative to other slices)

Under principle of difference...

$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{\sum_{\theta \in \mathcal{H}_0} P(x|\theta)P(\theta)} = \frac{P(x|\theta)}{P(x|\theta_1) + \dots + P(x|\theta_n)}$$

See this line... $P(\theta = .75 | x = 0.11)$

$$= \frac{0.191}{0.009 + 0.012 + 0.125 + 0.141 + 0.061} = \frac{.191}{.348} \approx 37\%$$

$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{MLE}}$ but $0.75 \neq 0.66$ Why? Poor choice of

prior did not cover all of the parameter space! $\mathcal{H}_0 \neq \mathcal{H} = (0,1)$
(for $Z = \text{bullet}$)

Main skeptic of Bayesian stats: Prior could be wrong!

Let's look at data one at a time $\Theta = \{0.25, 0.75\}$ (independent) $X_1 = 0$ (3)

What we know from ... prior and X_1

let $P(\Theta | X_1)$ be our new prior.

\Rightarrow No longer indifferent!

$X_2 = 1$

$$P(\Theta = 0.25 | X_1 = 0) = \frac{P(X_1 = 0 | \Theta = 0.25)}{P(X_1 = 0 | \Theta = 0.25) + P(X_1 = 0 | \Theta = 0.75)}$$

$$= \frac{0.25}{0.25 + 0.25} = 0.5$$

$$\Rightarrow P(\Theta = 0.75 | X_1 = 0) = 0.25 = 1 - P(\Theta = 0.25 | X_1 = 0)$$

$$P(\Theta = 0.25 | X_2 = 1) = \frac{P(X_2 = 1 | \Theta = 0.25) P(\Theta = 0.25 | X_1)}{P(X_2 = 1 | \Theta = 0.25) P(\Theta = 0.25 | X_1) + P(X_2 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_1)}$$

$$= \frac{0.25 \cdot 0.5}{0.25 \cdot 0.5 + 0.75 \cdot 0.25} = 0.5$$

we're back to square 1. For this prior, no information learned. Make sense?

Now we know prior, X_1, X_2 . Use this as prior

let $P(\Theta | X_2, X_1)$ be new prior when $X_3 = 1$.

$$P(\Theta = 0.25 | X_3 = 1) = \frac{P(X_3 = 1 | \Theta = 0.25) P(\Theta = 0.25 | X_2, X_1)}{P(X_3 = 1 | \Theta = 0.25) P(\Theta = 0.25 | X_2, X_1) + P(X_3 = 1 | \Theta = 0.75) P(\Theta = 0.75 | X_2, X_1)}$$

$$= \frac{0.25 \cdot 0.5}{0.25 \cdot 0.5 + 0.75 \cdot 0.5} = 0.25$$

Same as $P(\Theta = 0.25 | X = (0, 1, 1))$ from previously.

Is this true in general?

$$P(\Theta | X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n | \Theta) P(\Theta)}{P(X_1, \dots, X_n)}$$

Why? iid

$$= \frac{P(X_1 | \Theta) \dots P(X_n | \Theta) P(\Theta)}{P(X_1, \dots, X_n | X_1) P(X_1)} = P(\Theta | X_1)$$

$$= \frac{P(X_1 | \Theta) \dots P(X_n | \Theta) P(\Theta)}{P(X_1, \dots, X_n | X_1, X_2) P(X_1, X_2)} = P(\Theta | X_1, X_2)$$

etc...

New question: we have seen X_4 yet. Why is it here?

Of course $P(X_4 | \theta) = \theta^{x_4} (1-\theta)^{1-x_4}$ but... you don't know θ .

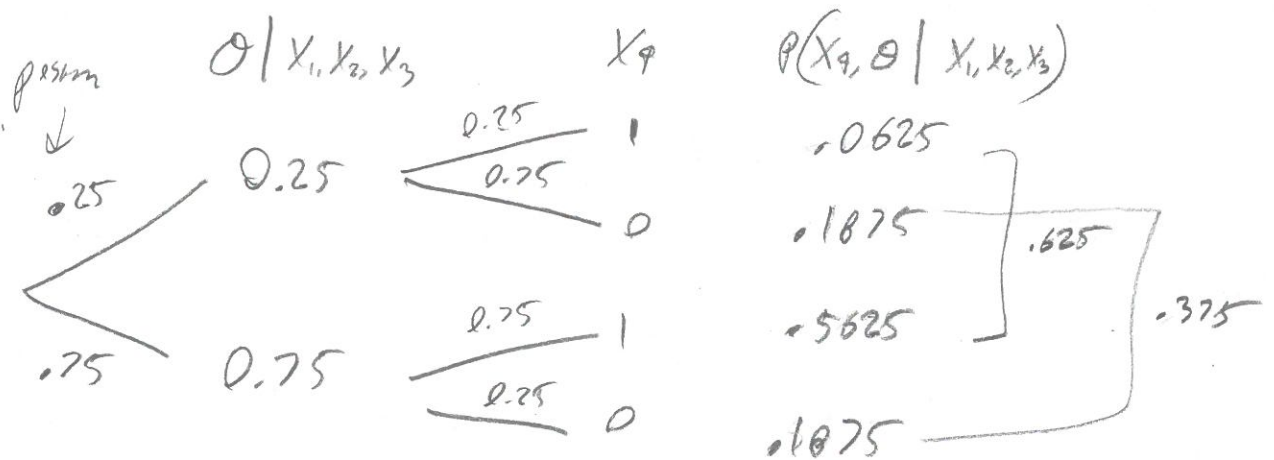
Previously... what did you do?

$$P(X_4 | \theta) \approx P(X_4 | \theta = \hat{\theta}_{MLE} = 0.66) = \text{Bern}(0.66)$$

What's the problem? Uncertainty in $\hat{\theta}_{MLE}$...

Bayesian Solution:

Seek: $P(X_4 | X_1, X_2, X_3)$



$$\Rightarrow P(X_4 | X_1, X_2, X_3) = \text{Bern}(0.625)$$

This incorporates all uncertainty of θ arising de prior & data.

What did we do?

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

See 95

$$P(Y) = \sum_{x \in \mathcal{X}} P(Y, x) \quad \text{but in the world of } X_1, X_2, X_3, \theta$$

Bayes Rule

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

why?

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

procedure: draw a θ from posterior, obtain X_4 under θ , repeat for all θ 's

Question rate on this procedure: draw θ from prior. Estimate how likely x_4 data is under θ relative to all possible θ 's, $P(X_1, X_2, X_3)$

for use this θ to see what X_4 is

$$= \sum_{\theta \in \Theta} P(X_4 | \theta) P(\theta)$$

"Posterior Predictive Distribution"

$$\begin{aligned}
 P(X_4 | \theta) & \stackrel{?}{=} P(X_4 | \theta, X_1, X_2, X_3) \\
 &= \frac{P(X_4, X_1, X_2, X_3, \theta)}{P(\theta, X_1, X_2, X_3)} \\
 &= \frac{P(X_1, X_2, X_3, X_4 | \theta) \cancel{P(\theta)}}{P(X_1, X_2, X_3 | \theta) \cancel{P(\theta)}} \quad \text{iid} \\
 &= \frac{P(X_1 | \theta) P(X_2 | \theta) P(X_3 | \theta) P(X_4 | \theta)}{\cancel{P(X_1 | \theta)} \cancel{P(X_2 | \theta)} \cancel{P(X_3 | \theta)}} \quad \checkmark
 \end{aligned}$$

Intuition: Once θ is known... data is useless when considering distr. of a future obs.

Generally...

$$P(X^* | X_1, \dots, X_n) = \sum_{\theta \in \Theta_0} P(X^* | \theta) P(\theta | X_1, \dots, X_n) \quad (\text{discrete})$$

$$\int_{\Theta_0} P(X^* | \theta) P(\theta | X_1, \dots, X_n) d\theta \quad (\text{continuous})$$

$\neq P(X^* | \hat{\theta}) \leftarrow$ Bad idea since one pt. estimate cannot represent the entire distr!

$$\hat{\theta}_{MLE} \neq \hat{\theta}_{MAP}$$

$$0.75 \neq 0.667 \quad \text{why?} \quad (H_0) \neq (H)$$

This is likely a bad idea... you shouldn't put 0 prob on elements in the param space without good reason!

So what prior can we use?

$$\theta \sim U(0,1) \quad \text{this has } \bullet \text{ Supp}(\theta) = [0,1] = (H)$$

like property!

• Principle of indifference... no θ 's are given special priority..

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

density function of std. unif.

$$x = (0, 1, 1)$$

$$P(\theta | x) = \frac{P(x_1, x_2, x_3 | \theta) P(\theta)}{P(x_1, x_2, x_3)} \propto P(x_1, x_2, x_3 | \theta) = \theta^2 (1-\theta)$$

$$\hat{\theta}_{MLE} = \arg \max \{ \theta^2 (1-\theta) \} \quad \text{get} \quad \theta = \frac{1}{\frac{d}{d\theta} [\theta^2 (1-\theta)]} = \frac{1}{\frac{d}{d\theta} [\theta^2 - \theta^3]} = \frac{1}{2\theta - 3\theta^2}$$

$$\Rightarrow 0 = 2 - 3\theta \Rightarrow \hat{\theta}_{MLE} = \frac{2}{3} \quad \checkmark$$

the "best model"

New question... what if I'm interested in $P(\theta \in [0.6, 0.7] | x)$ i.e. what we could've got before!

We need to solve for $P(\theta|x)$ explicitly... we just find its mode ...

$$P(\theta|x=(0,1,1)) = \frac{P(x_1, x_2, x_3|\theta) P(\theta)}{P(x_1, x_2, x_3)} = \frac{\theta^2(1-\theta)}{\int_0^1 P(x_1, x_2, x_3|\theta) P(\theta) d\theta} = 12\theta^2(1-\theta)$$

(H₀)
 $\sum P(x_1, x_2, x_3|\theta)$ integrating θ !

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

$$P(\theta \in [0.6, 0.7] | x) = \int_{0.6}^{0.7} P(\theta|x) d\theta = \int_{0.6}^{0.7} 12\theta^2(1-\theta) d\theta = 12 \left[\frac{\theta^3}{3} - \frac{\theta^4}{4} \right]_{0.6}^{0.7} = 12(0.0543 - 0.0396) = \boxed{0.1765}$$

this is the prob that true θ is between 0.6 and 0.7 assuming the prior.

Let's solve for general data x_1, \dots, x_n and same prior.

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)} = \frac{\prod_{i=1}^n P(x_i|\theta)}{\int_0^1 \prod_{i=1}^n P(x_i|\theta) d\theta} = \frac{\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}}{\int_0^1 \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} d\theta}$$

$$\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum x_i} (1-\theta)^{n - \sum x_i} \Rightarrow \text{prior only is dependent on } \sum x_i \Rightarrow \sum x_i \text{ from } n \text{ trials}$$

$$= \frac{\theta^{\sum x_i} (1-\theta)^{n - \sum x_i}}{\int_0^1 \theta^{\sum x_i} (1-\theta)^{n - \sum x_i} d\theta} \xrightarrow{\text{Beta function}} \text{Beta function } B(x, y) := \int_0^1 t^{x-1} (1-t)^{y-1} dt$$

↑
"conjugate β "