

# Lecture 15 PMH 390.03-a 4/4/16

## Semi Conjugacy

Situation

L1

You know  
there's no  
dependence

$\theta \sim N(150, 10^2)$ ,  $\sigma^2$  unknown so...

$\sigma^2 \sim \text{Jeffreys}$

$\sigma^2 \perp \theta$  a priori  $P(\theta, \sigma^2) = P(\theta)P(\sigma^2)$

when  $P(\theta)$  is conj prior for  $P(\sigma^2 | X, \theta)$   
and  $P(\sigma^2)$  is conj prior for  $P(\theta | X, \sigma^2)$

but  $P(\theta, \sigma^2)$  is not conj prior for  $P(\theta, \sigma^2 | X)$

$$\theta \sim N(\mu_0, \tau^2) \quad \sigma^2 \sim \text{InverseGamma}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})$$

$$P(\theta, \sigma^2 | X) \propto P(X | \theta, \sigma^2) P(\theta) P(\sigma^2)$$

$$\propto (\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2}((n-1)s^2 + (\bar{X} - \theta)^2)} e^{-\frac{1}{2\tau^2}(\theta - \mu_0)^2} (\sigma^2)^{-\left(\frac{\nu_0}{2} + 1\right)} e^{-\frac{\nu_0 \sigma_0^2}{2\sigma^2}}$$

$$\propto e^{-\frac{n}{2\sigma^2}(\bar{X} - \theta)^2} e^{-\frac{\theta^2}{2\tau^2}} e^{\frac{\theta \mu_0}{\tau^2}} (\sigma^2)^{-\left(\frac{\nu_0}{2} + 1\right)} e^{-\frac{\nu_0 \sigma_0^2 + (n-1)s^2}{2\sigma^2}}$$

$$= e^{-\frac{n}{2\sigma^2}\bar{X}^2} e^{\frac{n\bar{X}\theta}{\sigma^2}} e^{-\frac{\theta^2}{2\tau^2}} \dots$$

$$= e^{\left(\frac{n\bar{X}}{\sigma^2} + \frac{\mu_0}{\tau^2}\right)\theta - \left(\frac{1}{2\tau^2} + \frac{n}{2\sigma^2}\right)\theta^2} (\sigma^2)^{-\left(\frac{\nu_0}{2} + 1\right)} e^{-\frac{\nu_0 \sigma_0^2 + (n-1)s^2 + n\bar{X}^2}{2\sigma^2}}$$

$$-\frac{1}{2\tau^2} = a \quad d = \mu_0$$

$$\Rightarrow v = -\frac{1}{\tau^2} = \frac{1}{\tau^2 + \frac{n}{\sigma^2}}$$

$$-\frac{1}{2v}(\theta - d)^2 = -\frac{1}{2v}(\theta^2 - 2d\theta + d^2) = \frac{\theta^2}{-2v} + d\theta - \frac{d^2}{2v}$$

$$= P(\sigma^2 | X) \propto P(\theta | X)_{\text{non-sol.}}$$

$$\rightarrow \propto e^{-\frac{1}{2v}(\theta - d)^2} e^{\frac{d^2}{2v}}$$

$$P(\theta | X, \sigma^2) \rightarrow$$

$$\propto N\left(\theta, \sigma^2\right) e^{\frac{1}{2} \frac{\left(\frac{n\bar{X}}{\sigma^2} + \frac{\mu_0}{\tau^2}\right)^2}{\left(\frac{1}{\sigma^2} + \frac{n}{\sigma^2}\right)}} = \frac{(\tau^2 n \bar{X} + \sigma^2 \mu_0)^2}{(\sigma^2 \tau^2)^2} \frac{\sigma^2}{\sigma^2 + \tau^2}$$

$$\neq e^{-\frac{1}{2\sigma^2}}$$

$\Rightarrow$  not an IV, gamma

$\Rightarrow$  what do you do?

this derivation is mostly right... but may be off by a bit

$$\begin{array}{c} \text{HW} \\ \downarrow \\ P(\theta, \sigma^2 | x) \propto N(\theta, \sigma^2) \underbrace{k(\sigma^2 | x)} \end{array}$$

not a distribution ... can't sample from it

How about the following idea... Finite  $n$  approx. No  $n \rightarrow \infty$  req'd!

$$P(\theta | x) = c(x) k(\theta | x)$$

$$\ln P(\theta | x) = \ln c(x) + \underbrace{\ln(k(\theta | x))}_{g(\theta)}$$

Recall  $\text{Taylor}$   $\text{series}$   $\text{approx}$

$$f(x) \approx \text{Tay}(x, c, d) := \sum_{i=0}^d \frac{f^{(i)}(c)}{i!} (x-c)^i$$

$$g(\theta | x) \approx \text{Tay}(\theta, c, 2) = g(c | x) + g'(c | x)(\theta - c) + \frac{g''(c | x)(\theta - c)^2}{2}$$

What should  $c$  be?

Recall  $\hat{\theta}_{\text{MAP}} := \arg\max \{P(\theta | x)\} = \arg\max \{\ln P(\theta | x)\} = \arg\max \{\underbrace{\ln k(\theta | x)}_{g(\theta | x)}\}$

And  $g'(\hat{\theta}_{\text{MAP}} | x) = 0$

let  $c = \hat{\theta}_{\text{MAP}}$

$$\Rightarrow g(\theta | x) \approx g(\hat{\theta}_{\text{MAP}} | x) + g'(\hat{\theta}_{\text{MAP}} | x)(\theta - \hat{\theta}_{\text{MAP}}) + \frac{1}{2} g''(\hat{\theta}_{\text{MAP}} | x)(\theta - \hat{\theta}_{\text{MAP}})^2$$

$$\Rightarrow \ln P(\theta | x) \approx \ln c(x) + g(\hat{\theta}_{\text{MAP}} | x) + \frac{1}{2} g''(\hat{\theta}_{\text{MAP}} | x)(\theta - \hat{\theta}_{\text{MAP}})^2$$

$$\Rightarrow P(\theta | x) \approx c(x) e^{g(\hat{\theta}_{\text{MAP}} | x)} e^{\frac{1}{2} g''(\hat{\theta}_{\text{MAP}} | x)(\theta - \hat{\theta}_{\text{MAP}})^2}$$

$$\propto N\left(\hat{\theta}_{\text{MAP}}, \left(\frac{1}{g''(\hat{\theta}_{\text{MAP}} | x)}\right)^2\right)$$

Now you can sample from  $k(\theta | x)$ !

Does this make sense?

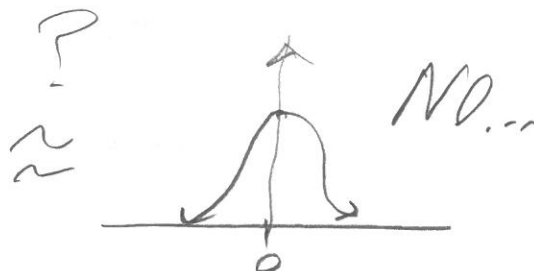
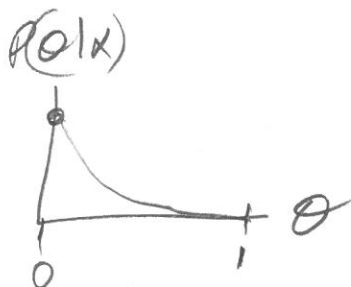
$$E[\theta|x] \approx \theta_{MAP} \text{ yes...}$$

Recall

$$\theta \sim \text{Beta}(1, 1)$$

$$x=0, n=3$$

$$\theta|x \sim \text{Beta}(1, 4)$$



Unless you really think the shape of the posterior is normal...  
all bets are off!!

Next idea

$$p(\theta, \sigma^2|x) \propto N(\theta_p, \sigma_p^2) k(\sigma^2|x)$$

Let's do the following... create a grid  $\{\theta_1, \dots, \theta_G\}$

every space

$$1) \text{ And calculate } \hat{c} := \left( \sum_{g=1}^G k(\theta_g|x) \right)^{-1}$$

then use...

$$p(\theta_g|x) \approx \hat{c} k(\theta_g|x)$$

In practice... keep all  $k(\theta_g|x)$  so we can  $\hat{F}(\theta_i|x) = \sum_{g=1}^G \hat{c} k(\theta_g|x)$

Draw  $u \sim U(0,1)$  and return  $\hat{F}^{-1}(u|x)$

## Disadvantages

① Numerically unstable..  $k(\theta|x) = 0$  or  $\infty$  in a computer

Sol: sample  $\ln k(\theta|x)$  and exponentiate afterwards

② If  $\text{Supp}(\theta)$  is an unbounded set,  
what support subset should we use?

$$\sigma_g^2 \in [0.0001, 1,000,000] ?$$

↑                      ↑  
need to make these  
decisions

If you make a wrong decision... you may miss some  
of the effective support.

(In our case, we know  $\sigma^2|x$  is unimodal, so  
we can end the grid when  $k(\sigma_g^2|x)$  gets small as  $\sigma_g^2 \uparrow$   
and we can start close to 0.)

③ In multiple dimensions, infeasible. Even  $1000^{10}$  is  $\infty$  in a computer.

For now: grid sampling okay... but we will reach  
this idea soon! We will take a break from Bayesian  
stuff to discuss the linear model: backbone of all stat. modeling.

Consider data from a bivariate distribution  $X, Y$

(5)

$\langle x_1, y_1 \rangle$   
 $\langle x_2, y_2 \rangle$   
 $\vdots$   
 $\langle x_n, y_n \rangle$

WLOG, let  $y$  be the  
"response", "outcome" or "dependent" variable  
and  $x$  is the  
"feature", "covariate", "regressor", "independent" variable

$X \xrightarrow[\text{change}]{\text{affects}} Y$  which may or may not be  
a causal effect

Causality beyond scope of  
course

$X \xrightarrow{f, \varepsilon} Y$  there is some function  $f$  and  
some noise generation  $\varepsilon \sim \mathcal{E}(\text{epsta})$   
 $Y = f(X) + \varepsilon$   
The goal is to model  $f$ .

$f(x) \in \mathcal{F}$  there is by game over stage!

But for our purposes now... restrict

$$f(x) \in \mathcal{F}_{\text{lin}} := \{ \beta_0 + \beta_1 x : \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R} \}$$

And restrict

$\varepsilon \neq h(x)$ , noise is independent

Before we get into random variables...