

Lecture 15

Let X_1, \dots, X_n ^{exch} $N(\theta, \sigma^2)$

We were looking at the posterior $P(\theta, \sigma^2 | x)$.

We had

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

under the informative prior, $p(\theta, \sigma^2) \propto N(\mu_0, \frac{\sigma^2}{m}) \text{InvGamma}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})$

and our posterior was

$$p(\theta, \sigma^2 | x) \propto N\left(\frac{n}{n+m} \bar{x} + \frac{m}{n+m} \mu_0, \frac{\sigma^2}{n+m}\right) \text{InvGamma}\left(\frac{\nu_0+n}{2}, \frac{\nu_0 \sigma_0^2 + (n-1)s^2 + \frac{m}{n+m}(\bar{x} - \mu_0)^2}{2}\right)$$

Thus we have

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

$$\text{NIG} \propto \text{Normal Likelihood} \cdot \text{NIG}$$

★ So the NIG is the conjugate prior for the normal likelihood ★

Now in the above case, $p(\theta, \sigma^2) = \underbrace{p(\theta | \sigma^2)}_{N(\mu_0, \frac{\sigma^2}{m})} \underbrace{p(\sigma^2)}_{\text{InvGamma}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})}$

so the distribution of θ depends on σ^2 .

What happens if $p(\theta, \sigma^2) = p(\theta) p(\sigma^2)$. I.e. if θ & σ^2 are independent?

$$\text{Let } p(\theta, \sigma^2) = \underbrace{p(\theta)}_{N(\mu_0, \tau^2)} \underbrace{p(\sigma^2)}_{\text{InvGamma}(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2})}$$

$$= N(\mu_0, \tau^2) \text{InvGamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$

where $\tau^2 \neq f(\sigma^2)$.

Then in this case we have

$$\begin{aligned} p(\theta, \sigma^2 | x) &\propto p(x | \theta, \sigma^2) p(\theta, \sigma^2) \\ &= p(x | \theta, \sigma^2) p(\theta) p(\sigma^2) \\ &= \text{NORMAL LIX} \cdot N(\mu_0, \tau^2) \cdot \text{InvGamma}\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right) \\ &\propto \underbrace{N(\theta_p, \sigma_p^2)}_{\text{stuff}} \cdot \text{stuff} \\ &\propto p(\theta | x, \sigma^2) \cdot k(\sigma^2 | x) \end{aligned}$$

$$\begin{aligned} \text{But } k(\sigma^2 | x) &\propto \text{InvGamma} \text{ \& } \\ k(\sigma^2 | x) &\propto \text{any known distribution} \end{aligned}$$

Now $p(\theta, \sigma^2 | x) \propto N(\theta_p, \sigma_p^2) \cdot \text{UNKNOWN DISTRIBUTION}$.

This leads to problems! How can we sample from it?

It's not a distribution we know so we can't sample from it.

We don't know what our posterior is!

Now pretend you only have kernel. So you only have $k(\sigma^2|x)$. By definition of kernel,

$$p(\sigma^2|x) = c(x) \cdot k(\sigma^2|x)$$

Remember, we only have $k(\sigma^2|x)$. If we can estimate $c(x)$ then we can estimate $p(\sigma^2|x)$ & hence estimate our posterior $p(\theta, \sigma^2|x)$.

$$p(\sigma^2|x) = c(x) \cdot k(\sigma^2|x)$$

$$\Rightarrow \ln p(\sigma^2|x) = \ln c(x) + \ln k(\sigma^2|x)$$

Call $g(\sigma^2|x) := \ln k(\sigma^2|x)$ so

$$\ln p(\sigma^2|x) = \ln c(x) + g(\sigma^2|x)$$

Now,

$$g(\sigma^2|x) \approx \text{Tay}(\sigma^2, c, 2) = \text{Tay}(\sigma^2, c, 2)$$

$$= \sum_{i=0}^2 \frac{g^{(i)}(c)}{i!} \cdot (\sigma^2 - c)^i$$

$$= g(c|x) + g'(c|x)(\sigma^2 - c) + \frac{g''(c|x)(\sigma^2 - c)^2}{2}$$

So we get a function which does NOT RELY on σ^2 .

Instead this relies on c . For what value should we choose c to be? Let $c = \hat{\sigma}_{\text{MAP}}^2$. Now

$$\hat{\sigma}_{\text{MAP}}^2 := \arg\max \{p(\sigma^2|x)\} = \arg\max \{k(\sigma^2|x)\} = \arg\max \{\ln(k(\sigma^2|x))\}$$

So $g'(\hat{\sigma}_{\text{MAP}}^2|x) = 0$ b/c $\hat{\sigma}_{\text{MAP}}^2$ is the maximum. Hence for $c = \hat{\sigma}_{\text{MAP}}^2$,

$$g(\hat{\sigma}_{\text{MAP}}^2|x) + g'(\hat{\sigma}_{\text{MAP}}^2|x)(\sigma^2 - \hat{\sigma}_{\text{MAP}}^2) + \frac{g''(\hat{\sigma}_{\text{MAP}}^2|x)(\sigma^2 - \hat{\sigma}_{\text{MAP}}^2)^2}{2}$$

$$= g(\hat{\sigma}_{\text{MAP}}^2|x) + \frac{g''(\hat{\sigma}_{\text{MAP}}^2|x)(\sigma^2 - \hat{\sigma}_{\text{MAP}}^2)^2}{2}$$

$$= \text{Tay}(\sigma^2, \hat{\sigma}_{\text{MAP}}^2, 2) \approx g(\sigma^2|x)$$

Now $g(\sigma^2|x) \approx g(\hat{\sigma}_{\text{MAP}}^2|x) + \frac{g''(\hat{\sigma}_{\text{MAP}}^2|x)(\sigma^2 - \hat{\sigma}_{\text{MAP}}^2)^2}{2}$. Now,

$$p(\sigma^2|x) \propto k(\sigma^2|x) = e^{g(\sigma^2|x)} \approx e^{g(\hat{\sigma}_{\text{MAP}}^2|x)} e^{\frac{1}{2} g''(\hat{\sigma}_{\text{MAP}}^2|x) (\sigma^2 - \hat{\sigma}_{\text{MAP}}^2)^2}$$

$$\propto e^{\frac{1}{2} g''(\hat{\sigma}_{\text{MAP}}^2|x) (\sigma^2 - \hat{\sigma}_{\text{MAP}}^2)^2}$$

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

approximate

So we have the distribution of $\sigma^2|x$! YAYYYY

Recall a function $f(x)$ can be approximated by a Taylor series centered at c of degree d .

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

where $\text{Tay}(x, c, d)$

More generally, $k(z|x) \propto N\left(\hat{z}_{\text{map}}, \frac{1}{\sqrt{g''(\hat{z}_{\text{map}}|x)}}\right)^2$

where $g(z|x) = \ln k(z|x)$.

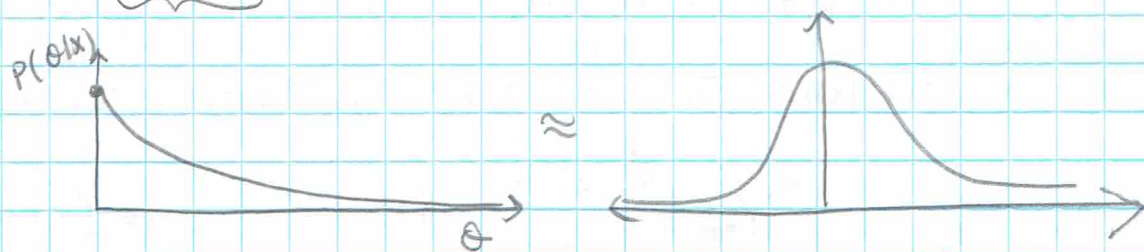
This approximation can be bad! Recall if

$X \sim \text{Bin}(n, \theta)$, $\theta \sim \text{Beta}(1, 1)$ then

$\theta|x \sim \text{Beta}(x+1, (n-x)+1)$. For $x=0, n=3$, we get

$\theta|x \sim \text{Beta}(1, 4)$. Note $\hat{\theta}_{\text{map}} = 0$

So $\theta|x \sim \text{Beta}(1, 4)$ & by above, $\theta|x \approx N(0, \text{something})$ so
the $\text{Beta}(1, 4)$ distribution $\approx N(0, \text{something})$



which is "false"

What I really mean is that the error is huge!!!

Now remember $P(\theta, \sigma^2|x) \propto \underbrace{p(\theta|x, \sigma^2)}_{\text{KNOWN}} \underbrace{k(\sigma^2|x)}_{\text{KNOWN what the function is but NOT know what the distribution is}}$

We want to draw/sample from $P(\theta, \sigma^2|x)$ but we need to know the distribution of $k(\sigma^2|x)$.

Here is the GRID METHOD.

(i.e. you don't know what the p.d.f. is)

① Let $D \subseteq \text{support of } \sigma^2$

② Fix $\epsilon > 0$

③ You know the function $k(\sigma^2|x)$

④ Define $\sigma_1^2 = \min(D)$
 $\sigma_2^2 = \min(D) + \epsilon$
 $\sigma_3^2 = \min(D) + 2\epsilon$
 \dots

$\sigma_n^2 = \max(D)$

⑤ You get $\{\sigma_1^2, \dots, \sigma_n^2\}$

⑥ Plugging into k yields

⑥ $P(\sigma_j^2|x) \approx \hat{c} k(\sigma_j^2|x)$

$\{k(\sigma_1|x), \dots, k(\sigma_n|x)\}$

where

$\hat{c} = \frac{1}{\sum_{j=1}^n k(\sigma_j|x)}$

⑦ $\hat{F}(\sigma_j^2|x) = \sum_{i=1}^j \hat{c} k(\sigma_i^2|x)$

⑧ $\sigma^2 = \hat{F}^{-1}(u)$

Problems with GRID APPROX

- ① If ϵ is too big
- ② If D is too small
- ③ If the kernel is multidimensional, then problems in estimating kernel.

Now let's discuss a time before probability, r.v.'s etc.

Consider a bivariate distribution X, Y

Let Y be the "outcome", "response", "dependent variable"

Let X be the "regressor", "feature", "covariate", "independent variable"

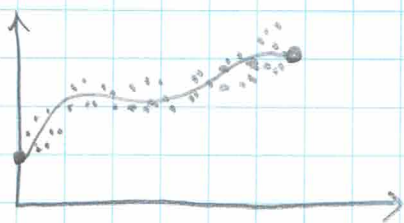
X affects Y

How does X affect Y ? By some function f
by some noise

Draws from a bivariate distribution look like

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

which is a 2-dimensional sample.



$$y = f(x) + \epsilon \quad \text{where } \epsilon \neq h(x)$$

primary goal: obtain $f(x)$
secondary goal: obtain ϵ

$$\text{Let } F = \{ f \mid f: \mathbb{R} \rightarrow \mathbb{R} \}$$

This set is too large, so limit the size of F .

$$F_0 = \{ \beta_0 + \beta_1 x : \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R} \}$$