

# MAT 4381

## Bayesian Inference

### Study Guide



### Winter 2024

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

## Bayesian vs. Frequentist

Frequentist: The parameter is not a random variable

Bayesian: The parameter is a random variable

## Bayes Theorem

$\Omega$  be the sample space  
 $E_1, E_2, \dots$  sequence of exhaustive events

$A \subset \Omega$  another event

$$P(E_k|A) = \frac{P(A|E_k)P(E_k)}{\sum_i P(A|E_i)P(E_i)}$$

## Bernoulli Trials

$X$  takes value 1 if event occurs and 0 if event do not occur

Bernoulli Distribution:  $Ber(\theta)$

$$P(X=1) = \theta$$

$$P(X=0) = 1-\theta$$

$$P(X=x) = f(x|\theta) = \theta^x(1-\theta)^{1-x} I(x \in \{0,1\}, \theta)$$

$$E[X] = \sum_{x=0}^1 x f(x|\theta) = \theta$$

$$\text{Var}(X) = \theta(1-\theta)$$

## Binomial Distribution

$$X_1, \dots, X_n \stackrel{iid}{\sim} Ber(\theta)$$

$$Y = \sum_{i=1}^n X_i \in \{0, 1, 2, \dots, n\} \sim \text{Bin}(n, \theta)$$

$$P(Y=y) = \binom{n}{y} \theta^y (1-\theta)^{n-y} I(y \in \{0, 1, \dots, n\})$$

$$E[Y] = n\theta \quad \text{Var}(Y) = n\theta(1-\theta)$$

## Bayes Formula for Distributions

$$f(\theta|y) = \frac{f(y|\theta)g(\theta)}{f_y(y)}$$

If  $\theta$  is discrete with p.m.f.  $g(\theta)$

$$f(\theta|y) = \frac{f(y|\theta)g(\theta)}{f_y(y)} = \frac{f(y|\theta)g(\theta)}{\sum_j f(y|\theta_j)g(\theta_j)}$$

If  $\theta \sim g(\theta)$  is continuous

$$f(\theta|y) = \frac{f(y|\theta)g(\theta)}{f_y(y)} = \int f(y|\theta)g(\theta)d\theta$$

The denominator is sometimes called the proportionality constant

$$f(\theta|y) = C f(y|\theta) g(\theta)$$

or  $f(\theta|y) \propto f(y|\theta)g(\theta)$  = likelihood  $\times$  prior

C some constant and  $\propto$  means "proportional to"

$g(\theta)$ : prior distribution for  $\theta$

$f(\theta|y)$ : posterior distribution

## Conjugate Prior

- A conjugate prior is also a convenient prior, since it allows one to easily calculate the posterior distribution

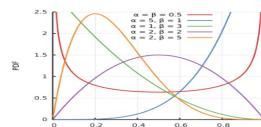
## Beta Distribution

$$g(\theta) = \frac{\Gamma(a+\beta)}{\Gamma(a)\Gamma(\beta)} \theta^{a-1} (1-\theta)^{\beta-1} I(0 < \theta < 1)$$

The distribution is the Beta( $a, \beta$ ) distribution

We use the fact that

$$\frac{\pi(a)}{\pi(a)} = \int_0^\infty x^{a-1} e^{-\beta x} dx, \quad a, \beta > 0$$



The Beta distribution is a conjugate prior for binomial distribution

If  $(Y|\theta=\theta) \sim f(y|\theta)$

$$f(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y} I(y \in \{0, 1, \dots, n\})$$

and  $(\theta | \alpha, \beta) \sim \text{Beta}(\alpha, \beta)$

we have

$$f(\theta|y) \propto f(y|\theta)g(\theta)$$

In other words,

$$f(\theta|y) \propto \theta^y (1-\theta)^{n-y} \theta^{a-1} (1-\theta)^{\beta-1}$$

The posterior is:

$$f(\theta|y) = \frac{\Gamma(a+\beta+n)}{\Gamma(a+y)\Gamma(\beta+n-y)} \theta^{a+y-1} (1-\theta)^{\beta+n-y-1}$$

If the prior is Beta( $\alpha, \beta$ ), then the posterior is Beta( $a+\alpha, \beta+n-y$ )

If  $\theta \sim \text{Beta}(\alpha, \beta)$ , then

$$E[\theta] = \frac{\alpha}{\alpha+\beta}$$

$$\text{Var}(\theta) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

$$\text{Mode: } \hat{\theta} = \frac{\alpha-1}{\alpha+\beta-2}$$

i) For large values of  $\alpha$  and  $\beta$  mode and mean are roughly the same.

ii) Median for Beta( $\alpha, \beta$ ) cannot be written in a closed form

iii) Peak of mode of beta density gets sharper as  $\alpha+\beta$  increases

iv) We need to use numerical methods to calculate the median for different values of  $\alpha$  and  $\beta$

## Credible Region

- We want to find I s.t.  $P(\theta \in I) = 1-\alpha$ , where  $1-\alpha$  is the confidence level
- We can also use simulation method

## Choosing Beta function

First, consider priors that don't express a strong opinion.

Later, consider priors which incorporate a strong prior belief. If there are no preferences on the value for  $\theta$ , then one might consider using a

"flat prior". That is, the uniform distribution between 0 and 1 ( $Beta(1,1)$ )

In some sense, this implies that each value of  $\theta$  is equally likely

Since the prior is flat, the posterior distribution has the same shape as the likelihood multiplied by a constant C

$$h(\theta|x) = C f(x|\theta)g(\theta) \propto f(x|\theta)$$

Since the uniform distribution is Beta( $1,1$ ), then the posterior is Beta( $y+1, n-y+1$ )

## Making Bayesian Inference

The posterior contains the results of the analysis

To communicate the result, we can either present the entire posterior distribution, or give summary statistics of the posterior

For credible regions, there are 2 common types:

- Smallest interval
- Equal tail area

## Jeffreys Prior

Jeffreys proposed that an acceptable "non-informative prior finding principle" should be invariant under monotone transformations of the parameter

The Jeffreys prior satisfying this invariance is proportional to  $\sqrt{I(\theta)}$ , where  $I(\theta)$  is the Fisher information on  $\theta$ . Can be generalized when  $\theta$  is k-dimensional to  $\sqrt{\det(I(\theta))}$

$$I(\theta) = -E\left[\frac{\partial^2 \ln f(x|\theta)}{\partial \theta^2}\right]$$

For k-dimensional, we have

$$I(\theta) = -E\left[\left(\frac{\partial^2 \ln f(x|\theta)}{\partial \theta_i \partial \theta_j}\right)_{(i,j)}\right]$$

## Multinomial Distribution

$$(X_1, \dots, X_k) \sim \text{Multinomial}(n, P_1, \dots, P_k)$$

$$P_k = P_1 + P_2 + \dots + P_{k-1}$$

Models probability of counts for each side of an experiment with k exhaustive outcomes

$$f(X_1, X_2, \dots, X_k) = \binom{n}{X_1, X_2, \dots, X_k} P_1^{X_1} P_2^{X_2} \dots P_k^{X_k} (1-P_1-P_2-\dots-P_{k-1})^{X_k}$$

where  $X_i = 0, 1, 2, \dots, n$ ,  $X_1 + X_2 + \dots + X_k = n$

## Simple Case: Trinomial Distribution

$$(X, Y) \sim \text{Multinomial}(n, P_1, P_2)$$

$$f(x,y) = \binom{n}{x,y} P_1^x P_2^y (1-P_1-P_2)^{n-x-y}$$

$$X, Y = 0, 1, 2, \dots, n, X+Y \leq n$$

We can also get that

$$X \sim \text{Bin}(n, P_1)$$

$$Y \sim \text{Bin}(n, P_2)$$

$$Y|X=x \sim \text{Bin}(n-x, \frac{P_2}{1-P_1})$$

$$E[Y|X=x] = \frac{P_2}{1-P_1} (n-x)$$

$$\text{Var}[Y|X=x] = \frac{P_2}{1-P_1} (1-P_2)(n-x)$$

$$E[X] = \frac{P_1}{1-P_1} n, E[Y] = \frac{P_2}{1-P_1} n$$

$$\text{Cov}(X, Y) = -n(n+1) P_1 P_2$$

$$P(X, Y) = \frac{-n(n+1) P_1 P_2}{\sqrt{P_1(1-P_1)P_2(1-P_2)}} = \frac{-(n+1) P_1 P_2}{\sqrt{P_1 P_2 (1-P_1) (1-P_2)}}$$

## Jeffreys Prior for Trinomial Distribution

$$g(P_1, P_2) \propto \frac{1}{\sqrt{P_1 P_2 (1-P_1) (1-P_2)}}$$

Special case of Dirichlet distribution  
 $\text{Dir}(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$

## Dirichlet Distribution

$$(X_1, X_2, \dots, X_k) \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_k)$$

$$f(x_1, \dots, x_k) = \frac{\prod_{i=1}^k x_i^{\alpha_i}}{\prod_{i=1}^k \Gamma(\alpha_i)} \left(\prod_{i=1}^k x_i\right)^{1-\sum_{i=1}^k \alpha_i}$$

Defined on  $S = \{(x_1, \dots, x_k) : \sum_{i=1}^k x_i \leq 1, x_i \geq 0, i=1, \dots, k\}$

## Posterior for Trinomial

$$(\theta_1, \dots, \theta_{k-1} | X_1=x_1, \dots, X_k=x_{k-1}) \sim \text{Dir}(\alpha_1+x_1, \dots, \alpha_{k-1}+x_{k-1}, n-\sum x_i)$$

$$E[\theta_i | X_1=x_1, \dots, X_{k-1}=x_{k-1}] = \frac{\alpha_i + x_i}{\sum_{j=1}^k \alpha_j + n}$$

$$\theta = [\theta_1, \theta_2, \dots, \theta_{k-1}, \theta_k]^T$$

$$\theta_k = 1 - \theta_1 - \dots - \theta_{k-1} \quad x_k = n - x_1 - \dots - x_{k-1}$$

$$\text{Notice: } (\theta_1, \theta_2) \sim \text{Dir}(\alpha_1, \alpha_2, \alpha_3 + \dots + \alpha_k)$$

## Basic Properties of Dirichlet Distribution

$$(Y_1, \dots, Y_{k-1}) \sim \text{Dir}(\alpha_1, \dots, \alpha_{k-1}) \quad \alpha_0 = \sum_{i=1}^k \alpha_i$$

$$E[Y_i] = \frac{\alpha_i}{\alpha_0} \quad \text{Var}[Y_i] = \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$$

$$\text{Cov}(X_i, X_j) = -\frac{\alpha_i \alpha_j}{\alpha_0^2(\alpha_0 + 1)}$$

## Conjugate Prior and Posterior for Poisson

$$X_1, \dots, X_n \stackrel{iid}{\sim} f(x|\theta) = \frac{\exp(-\theta)x^x}{x!} I(x=0, 1, \dots)$$

$$\text{Likelihood: } \prod_{i=1}^n f(x_i|\theta) = \prod_{i=1}^n \frac{\exp(-\theta)x_i^x}{x_i!} \propto \exp(-n\theta) \theta^y$$

$$y = \sum_{i=1}^n x_i$$

$$\text{Conjugate Prior: } g(\theta) = \frac{\theta^x}{\Gamma(x)} \theta^{x-1} \exp(-\theta) I(\theta > 0)$$

$$\text{Posterior: } h(\theta|y) \propto \theta^{y+x+1} \exp(-(\theta+n)\theta) I(\theta > 0)$$

$$h(\theta|y) = \frac{(\theta+n)^{y+x}}{\Gamma(y+x)} \theta^{y+x+1} \exp(-(\theta+n)\theta) I(\theta > 0)$$

$$E[\theta|Y=y] = \frac{y+x}{n+p}$$

$$\text{Var}[\theta|Y=y] = \frac{y+x}{(n+p)^2}$$

$$\text{Mode}(\theta|Y=y) = \frac{y+x-1}{n+p}$$

## Jeffrey's Prior for Poisson

$$f(x|\lambda) = \frac{\exp(-\lambda)\lambda^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}$$

$$-\mathbb{E}\left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta)\right] = \frac{n}{\lambda}$$

Jeffrey's prior:  $g(\lambda) \propto \frac{1}{\lambda}$

This leads to an improper prior since we have

$$\int_0^\infty \frac{d\lambda}{\sqrt{\lambda}} = \infty$$

No normalizing constant will turn it into a prior p.d.f.

We can still use it to calculate the posterior

$$\lambda | X \propto \exp(-n\lambda) \lambda^{\sum_{i=1}^n x_i - 1}$$

The use of improper prior leads to Gamma posterior

$$\lambda | X \sim \frac{1}{\Gamma(\sum_{i=1}^n x_i + 1)} \lambda^{\sum_{i=1}^n x_i - 1} \exp(-n\lambda)$$

$$\mathbb{E}[\lambda | X_1, \dots, X_n] = \frac{\sum_{i=1}^n x_i + 0.5}{n}$$

$$\text{Var}[\lambda | X_1, \dots, X_n] = \frac{\sum_{i=1}^n x_i + 0.5}{n^2}$$

## Normal Distribution

$$(Y_1, \dots, Y_n | \mu) \sim f(y_1, \dots, y_n | \mu) \propto \exp\left(-\sum_{i=1}^n \frac{(y_i - \mu)^2}{2\sigma^2}\right)$$

$\sigma^2$  known

$$f(y_1, \dots, y_n | \mu) \propto \exp\left(-\frac{n(\mu - \bar{y})^2}{2\sigma^2}\right)$$

The Jeffrey's prior is

$$g(\mu) \propto 1$$

This constitutes an improper prior

However, we can still calculate the posterior

$$(\mu | Y_1 = y_1, \dots, Y_n = y_n) \sim N(\bar{y}, \frac{\sigma^2}{n})$$

## The normal prior

$$Y_1, \dots, Y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$\sigma^2$  known

$$\text{Prior: } g(\theta) \propto \exp(-(\theta - \mu)^2 / \gamma^2)$$

$$\text{Likelihood: } \prod_{i=1}^n f(y_i | \mu, \sigma^2) \propto \exp\left(-\frac{(y_i - \mu)^2}{2\sigma^2}\right)$$

Posterior for  $\theta$

$$\begin{aligned} h(\theta | y_1, \dots, y_n) &\propto \exp\left(-\frac{(\theta - \mu)^2}{\gamma^2}\right) \cdot \exp\left(-\frac{(y_1 - \mu)^2}{2\sigma^2}\right) \\ &\propto \exp\left(\frac{(\theta - \mu)^2}{\gamma^2} - \frac{(y_1 - \mu)^2}{2\sigma^2}\right) \\ &\propto \exp\left(-\frac{\theta^2 - 2\theta\mu + \mu^2}{\gamma^2} - \frac{y_1^2 - 2y_1\mu + \mu^2}{2\sigma^2}\right) \\ &\propto \exp\left(-\frac{\theta^2 - 2\theta\mu + \mu^2 + 2y_1^2 - 4y_1\mu + 2\mu^2}{\gamma^2 + 2\sigma^2}\right) \end{aligned}$$

$$(y_1, \dots, y_n) \sim N\left(\frac{\gamma^2\mu + 2\sigma^2y_1}{\gamma^2 + 2\sigma^2}, \frac{\sigma^2}{\gamma^2 + 2\sigma^2}\right)$$

$$\mathbb{E}[\theta | Y_1 = y_1, \dots, Y_n = y_n] = \frac{\gamma^2\mu + 2\sigma^2y_1}{\gamma^2 + 2\sigma^2}$$

$$\text{MAP} = \mathbb{E}[\theta | Y_1 = y_1, \dots, Y_n = y_n] = \text{Median}(\theta | Y_1 = y_1, \dots, Y_n = y_n)$$

$$\text{Posterior precision: } \frac{\tau^2 + \sigma^2/n}{\tau^2 + \sigma^2} = \frac{1}{\sigma^2/n} + \frac{1}{\tau^2}$$

Posterior precision is equal to prior precision plus observation precision

Since  $E[\theta | y_1, \dots, y_n] = \frac{\tau^2\mu + \sigma^2y_1}{\tau^2 + \sigma^2} + \frac{\tau^2}{\tau^2 + \sigma^2} \bar{y}$ , the posterior mean is the weighted average of the prior mean and the mean of observations, where the weights are the proportions of the precision to the posterior precision.

Credible region for  $\theta$ :

$$\frac{\tau^2\bar{y} + \sigma^2y_1/n}{\tau^2 + \sigma^2/n} \pm \frac{z_{\alpha/2}}{\sqrt{\tau^2 + \sigma^2/n}}$$

If  $\sigma^2$  unknown, we can use

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

We change  $Z_{\alpha/2}$  to  $t_{\alpha/2}(n-1)$

Credible region becomes

$$\frac{\tau^2\bar{y} + \sigma^2y_1/n}{\tau^2 + \hat{\sigma}^2/n} \pm t_{\alpha/2}(n-1) \frac{\sqrt{\hat{\sigma}^2/n}}{\sqrt{\tau^2 + \hat{\sigma}^2/n}}$$

## Jeffrey's prior for $\sigma$ and $\sigma^2$

Suppose  $\mu$  is fixed

Then we get the Jeffrey's prior for  $\sigma$  and  $\sigma^2$  are

$$g(\sigma) \propto \frac{1}{\sigma} \quad \text{and} \quad g(\sigma^2) \propto \frac{1}{\sigma^2}$$

## Bayes Factor and Relative Evidence

In model selection, say 2 models.

The Bayes factor is defined by

$$BF(1,2) = \frac{P(M_1 | \text{Data}) P(M_2)}{P(M_2 | \text{Data}) P(M_1)}$$

With a prior on model, the Bayes factor can be calculated. With equal probability on models, we get

$$BF(1,2) = \frac{P(M_1 | \text{Data})}{P(M_2 | \text{Data})}$$

## Inverse $\chi^2$ Distribution

$V \sim \chi^2(r)$ . The p.d.f. for  $V$  is

$$f(v) = \frac{1}{2^{r/2} \Gamma(r/2)} v^{r/2-1} I(v > 0)$$

$$\text{Define } U = \frac{1}{V}$$

We say  $U \sim \text{Inv-}\chi^2(v)$

The p.d.f. for inverse- $\chi^2$  is

$$g(u) = \frac{1}{2^{r/2} \Gamma(r/2)} u^{-r/2-1} I(u > 0)$$

$$E[U] = \frac{1}{r-2} \quad \text{Mode}(U) = \frac{1}{r+2}$$

$$\text{Var}(U) = \frac{2}{(r-2)^2(r-4)}$$

Mean exists if  $r > 2$  and  $r > 4$

## Inverse Gamma Distribution

$$X \sim \text{Gamma}(\alpha, \beta)$$

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} \exp(-\beta x) x^{\alpha-1} I(x > 0)$$

Define  $U = X^{-1}$ . The p.d.f. for  $U$  is

$$g(u) = \frac{\beta^\alpha}{\Gamma(\alpha)} \exp(-\beta/u) u^{-\alpha-1} I(u > 0)$$

Denote by  $U \sim \text{IG}(\alpha, \beta)$

$$E[U] = \frac{\beta}{\alpha-1} \quad \text{Mode}(U) = \frac{\beta}{\alpha+1}$$

$$\text{Var}(U) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}$$

## Inverse Gamma prior on $\sigma^2$ , $\mu$ fixed

$$(y_1, \dots, y_n | \mu, \sigma^2) \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$g(\sigma^2) = \frac{\beta^\alpha}{\Gamma(\alpha)} \exp(-\frac{\beta}{\sigma^2}) (\sigma^2)^{-\alpha-1}$$

$$h(\sigma^2 | \mu, y_1, \dots, y_n) \propto (\sigma^2)^{-\alpha-1} \exp\left(-\frac{\sum (y_i - \mu)^2}{2\sigma^2}\right) \cdot \exp(-\beta/\sigma^2) \cdot \sigma^{\alpha-1}$$

$$\propto (\sigma^2)^{-(\alpha+\frac{1}{2})} \exp\left(-\frac{\beta + \sum (y_i - \mu)^2}{2\sigma^2}\right)$$

Therefore

$$(\sigma^2 | \mu, y_1, \dots, y_n) \sim \text{IG}\left(\alpha + \frac{n}{2}, \beta + \frac{\sum (y_i - \mu)^2}{2}\right)$$

In this case

$$E[\sigma^2 | \mu, y_1, \dots, y_n] = \frac{\beta + \sum_{i=1}^n (y_i - \mu)^2 / 2}{\alpha + n/2 - 1}$$

$$\text{Mode} = \frac{\beta + \sum_{i=1}^n (y_i - \mu)^2 / 2}{\alpha + n/2 + 1}$$

If  $X \sim \chi^2(v)$ , then  $1/X \sim \text{Inv-}\chi^2(v)$

Inverse chi-squared distribution is a special case of gamma distribution with parameters  $\alpha = v/2$ ,  $\beta = 1/2$ , or  $\alpha = v/2$ ,  $\beta = v^2/2$  for scaled inverse- $\chi^2$  distribution

## Posterior Predictive Distribution

Let  $y_{n+1}$  be the next observations drawn after the random sample  $Y_n = (y_1, y_2, \dots, y_n)$ .

The predictive accuracy for  $(y_{n+1} | y_1, y_2, \dots, y_n)$  is the conditional density

$$f(y_{n+1} | Y_n) = \int f(y_{n+1} | \mu) g(\mu | Y_n) d\mu$$

## Jeffrey's Prior for normal family with $\mu$ and $\sigma^2$ unknown

$$X \sim N(\mu, \sigma^2) \quad \Theta = [\mu \ \sigma^2]^T$$

Then we get

$$-E\left[\frac{\partial^2 \ln f}{\partial \mu^2}\right] = \frac{1}{\sigma^2} \quad -E\left[\frac{\partial^2 \ln f}{\partial \sigma^2}\right] = 0$$

$$-E\left[\frac{\partial^2 \ln f}{\partial \mu \partial \sigma^2}\right] = \frac{1}{\sigma^4}$$

The Jeffrey's prior is

$$g(\mu, \sigma^2) \propto (\sigma^2)^{-\frac{3}{2}}$$

To calculate posterior, we get

$$h(\mu, \sigma^2 | y_1, \dots, y_n) \propto (\sigma^2)^{-\alpha-1} \exp\left(-\frac{(n-1)\sigma^2 + (\mu - \bar{y})^2}{2\sigma^2}\right)$$

This shows that

$$(\mu, \sigma^2 | y_1, \dots, y_n) \sim N - \text{IG}(\mu = \bar{y}, \lambda = n, \alpha = n, \beta = \sum_{i=1}^n (y_i - \bar{y})^2)$$

## Normal Inverse Gamma Distribution

Let  $\lambda > 0$  and

$$(y | \mu, \sigma^2, \lambda) \sim N(\mu, \sigma^2/\lambda)$$

$$(\sigma^2 | \mu, \lambda) \sim \text{IG}(\mu, \lambda)$$

Then we say  $(y, \sigma^2)$  has normal inverse gamma distribution with parameters  $(\mu, \lambda, \alpha, \beta)$  which is denoted by

$$(y, \sigma^2) \sim \text{NIG}(\mu, \lambda, \alpha, \beta)$$

$$(y, \sigma^2) \sim N - \text{IG}(\mu, \lambda, \alpha, \beta)$$

We can write the joint p.d.f. for  $(y, \sigma^2)$  as follows

$$f(y, \sigma^2 | \mu, \lambda, \alpha, \beta) = \frac{1}{\sqrt{2\pi}} \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp\left(-\frac{\beta + (y - \mu)^2}{2\sigma^2}\right)$$

year,  $\sigma^2 > 0$

By definition,  $(\sigma^2 | \mu, \lambda) \sim \text{IG}(\mu, \lambda)$

We can get that  $y \sim t(2\alpha, \mu, \beta/\lambda\alpha)$

## Student t distribution

Let  $W \sim N(0, 1)$  and  $V \sim \chi^2(r)$  and  $W$  and  $V$  are mutually independent. The joint distribution for  $W$  and  $V$  is

$$f(v, w) = \frac{1}{\sqrt{2\pi r}} \frac{1}{2^{r/2} \Gamma(r/2)} \exp\left(-\frac{w^2}{2}\right) V^{r/2-1} \exp(-v^2/2)$$

for  $w \in R$ ,  $v > 0$

If we define  $T = \frac{W}{\sqrt{V}}$  and  $U = V$ ,

we get

$$g(u, t) = \frac{1}{\sqrt{\pi t}} \frac{1}{2^{(r+1)/2}} u^{(r+1)/2-1} \exp(-u(u + \frac{t}{u}))$$

$$g(t) = \frac{1}{\sqrt{\pi t}} \frac{\Gamma((r+1)/2)}{\Gamma(r/2)(1 + \frac{t}{u})^{(r+1)/2}}$$

when  $t > 0$

This is known as the student t distribution

Notes:

- $g(t)$  is a Cauchy distribution if  $r=1$
- $g(t) \rightarrow N(0, 1)$  as  $r \rightarrow \infty$

If we use the transformation  $V = \sigma^2 t + \mu$  when  $\sigma > 0$ , we get

$$h(v) = \frac{1}{\sigma \sqrt{\pi r}} \frac{\Gamma((r+1)/2)}{\Gamma(r/2)} \frac{1}{(1 + \frac{v-\mu}{\sigma})^{(r+1)/2}}$$

Denote this distribution by  $t(r, \mu, \sigma)$

If we let  $U = \sigma \left( \frac{W}{\sqrt{V}} \right) + \mu$ , we get

$$\text{Median}(U) = \text{Mode}(U) = E[U] = \sigma E[W] \frac{1}{\sqrt{V}} + \mu$$

If  $E[\frac{1}{\sqrt{V}}]$  exists then  $E[U] = \mu$

$$\text{Var}(U) = \frac{r\sigma^2}{r-2}$$

## Properties of Normal Inverse Gamma Distribution

$$(y, \sigma^2) \sim N - \text{IG}(\mu, \lambda, \alpha, \beta) \Rightarrow (\sigma^2 | \mu, \lambda) \sim \text{IG}(\mu, \lambda)$$

$$E[\sigma^2 | \alpha, \beta] = \frac{\beta}{\alpha-1}$$

$$\text{Var}[\sigma^2 | \alpha, \beta] = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}$$

$$\text{Var}[y] = E\left[\frac{\sigma^2}{\lambda}\right] = \frac{\beta}{\lambda(\alpha-1)}$$

## Normal-Inverse Gamma Prior

$$(y_1, \dots, y_n | \mu, \sigma^2) \sim N(\mu, \sigma^2)$$

$$\mu | \sigma^2 \sim N(\mu_0, \sigma^2/n)$$

$$\sigma^2 \sim IG(\alpha, \beta)$$

$n$ : equivalent sample size of the prior

$n$ : equivalent sample size we would obtain for  $n$  observations

The posterior is

$$f(\mu, \sigma^2 | y_1, \dots, y_n) \propto (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{\sum_i (y_i - \mu)^2}{2\sigma^2}\right)$$

$$\times (\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{n(\mu - \mu_0)^2}{2\sigma^2}\right)$$

$$\times (\sigma^2)^{-(\alpha+\beta)} \exp\left(\frac{\beta}{\sigma^2}\right)$$

Therefore,

$$f(\mu, \sigma^2 | y_1, \dots, y_n) \propto (\sigma^2)^{-(n+2\alpha+3)/2}$$

$$\times \exp\left(-\frac{\sum_i (y_i - \mu)^2}{2\sigma^2} - \frac{n(\mu - \mu_0)^2}{2\sigma^2} - \frac{\beta}{\sigma^2}\right)$$

Expand to get

$$f(\mu, \sigma^2 | y_1, \dots, y_n) \propto \exp\left(-\mu^2 \left(\frac{(n+n_0)\mu}{2\sigma^2} + \frac{(n_0\mu_0\mu_0)}{\sigma^2}\right)\right)$$

$$\times (\sigma^2)^{-(n+2\alpha+3)/2} \exp\left(-\frac{\sum_i (y_i - \mu)^2 + n(\mu - \mu_0)^2 + 2\beta}{2\sigma^2}\right)$$

Therefore

$$f(\mu, \sigma^2 | y_1, \dots, y_n) \propto \exp\left(-\frac{\mu^2 - 2\mu \left(\frac{(n+n_0)\mu}{2\sigma^2} + \frac{(n_0\mu_0\mu_0)}{\sigma^2}\right)}{2\sigma^2}\right)$$

$$\times (\sigma^2)^{-(n+2\alpha+3)/2} \exp\left(-\frac{\sum_i (y_i - \mu)^2 + n(\mu - \mu_0)^2 + 2\beta}{2\sigma^2}\right)$$

In other words,

$$f(\mu, \sigma^2 | y_1, \dots, y_n) \propto \exp\left(-\frac{(\mu - \frac{(n+n_0)\mu}{2\sigma^2})^2}{2\sigma^2}\right)$$

$$\times (\sigma^2)^{-(n+2\alpha+3)/2} \exp\left(-\frac{\sum_i (y_i - \mu)^2 + n(\mu - \mu_0)^2 + 2\beta}{2\sigma^2} + \frac{(n_0\mu_0\mu_0)}{\sigma^2}\right)$$

The posterior for  $\sigma^2$  is

$$f(\sigma^2 | y_1, \dots, y_n) \propto (\sigma^2)^{-(n+2\alpha+3)/2} \exp\left(-\frac{\sum_i (y_i - \bar{y})^2 + n(\bar{y} - \mu)^2}{2\sigma^2} + \frac{(n_0\mu_0\mu_0)}{2\sigma^2} + \frac{(n\mu - n\mu_0)^2}{2\sigma^2}\right)$$

For the posterior of  $\mu$ , we have

$$h(\mu | \sigma^2, \bar{y}) = \frac{f(\mu, \sigma^2 | \bar{y})}{f(\sigma^2 | \bar{y})}$$

$$\mu | \sigma^2, \bar{y} \sim N\left(\frac{n\bar{y} + n_0\mu_0}{n+n_0}, \frac{\sigma^2}{n+n_0}\right)$$

$$h(\mu | \bar{y}) = \int_0^\infty h(\mu | \sigma^2, \bar{y}) h(\sigma^2 | \bar{y}) d\sigma^2$$

Simulating Posterior for  $\mu$  and  $\sigma^2$  given data

$$\sigma^2 | y_1, \dots, y_n \sim IG(\alpha + \frac{n}{2}, \beta + \frac{\sum_i (y_i - \bar{y})^2}{2} + \frac{n(n\bar{y} - n\mu_0)^2}{2\sigma^2})$$

$$\mu | \sigma^2, y_1, \dots, y_n \sim N\left(\frac{n\bar{y} + n_0\mu_0}{n+n_0}, \frac{\sigma^2}{n+n_0}\right)$$

We can first simulate  $\sigma^2$  from the inverse gamma distribution. Then we use it to simulate  $\mu$  from the normal distribution.

Algorithm:

- Load data  $y_1, \dots, y_n$
- Find  $\bar{y}$ ,  $\alpha$ , and  $\beta$  parameters
- Simulate  $\sigma^2$  a number of times
- Use  $\sigma^2$  to simulate  $\mu$

We have  $\mu$  is symmetrically distributed. We can find a credible region the following way:

- Sort the simulated  $\mu$  values
- Find 2.5 and 97.5 percentiles
- Then we get the credible region

This method is more tricky with  $\sigma^2$  since it is not symmetrically distributed.

## Acceptance-Rejection Sampling

If we only know the unscaled target distribution, we can use acceptance-rejection sampling to draw a random sample from the target distribution.

This method works by drawing a random sample from an easily sample candidate distribution ( $g(\theta)$ ) by accepting some which satisfy a criterion. Let  $M$  be a positive constant and  $g_0(\theta)$  be the easily sampled candidate and

$$g(\theta) f(y|\theta) \leq M g_0(\theta)$$

The goal is to draw a sample from the posterior

$$h(\theta | y) \propto g(\theta) f(y|\theta)$$

The algorithm is as follows:

- Draw a sample from  $g_0(\theta)$
- Calculate  $w(\theta) = \frac{g(\theta) f(y|\theta)}{M g_0(\theta)}$
- Note:  $w(\theta) \sim (0, 1)$
- Draw  $u \sim \text{Uniform}(0, 1)$
- If  $u < w(\theta)$ , accept  $\theta$
- Otherwise, reject  $\theta$  and iterate step 1 until the first acceptance.

Theorem: Let  $f$  and  $g$  be two p.d.f.s with c.d.f.s  $F$  and  $G$ , respectively.

Assume  $f(y) \leq Mg(y)$  for all  $y$ . Draw  $U \sim \text{Uniform}(0, 1)$  and  $Y \sim g(y)$  independent from  $U$ . If  $U < \frac{f(y)}{Mg(y)}$  then accept  $Y$ , otherwise, ignore and iterate the previous steps until a sample is accepted. Define  $N$  to be the number of repeats until the first sample acceptance. We have:

- $N \sim \text{Geom}(p)$
- $P(N=n) = p(1-p)^{n-1} I(n \in \{1, 2, \dots\})$
- with  $p = \frac{1}{M}$
- $P(Y \leq y | U < \frac{f(y)}{Mg(y)}) = F(y)$

### Plug in Principle

Let  $X_1, \dots, X_m$  be a sequence of i.i.d. random variables with unknown distribution  $F$ . A natural estimate for  $F$  is

$$F_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

The method that replaces the empirical process, a natural estimate, for the actual distribution, is usually called the plug in principle. In other words the parameter  $T(F)$  is estimated by the statistic  $T(F_n)$ . In general from the plug in principle,

$$T(F) = \int g(x) dF(x)$$

should be estimated by

$$T(F_n) = \int g(x) dF_n = \frac{g(x_1) + \dots + g(x_n)}{n}$$

This roots from the fact that  $F_n \xrightarrow{a.s.} F$  (strong law of large numbers)

### Alternative Methodology

The expert makes a prior assumption on  $F$  (i.e.  $F$  belongs to a family of distributions) and then it will be updated after observing data.

## Dirichlet Process

Consider a space  $\mathcal{X}$  with a  $\sigma$ -algebra  $A$  of subsets of  $\mathcal{X}$ . Let  $H$  be a fixed probability measure on  $(\mathcal{X}, A)$  and  $\alpha$  be a positive number.

A random probability measure  $P = \{P(A)\}_{A \in A}$  is called a DP( $\alpha H$ ), if for any finite measurable partition  $\{A_1, \dots, A_n\}$  of  $\mathcal{X}$

$$(P(A_1), \dots, P(A_n)) \sim \text{Dir}(\alpha H(A_1), \dots, \alpha H(A_n))$$

### $P \sim DP(\alpha, H)$

For any  $A \in A$ , we have

$$P(A) \sim \text{Beta}(\alpha H(A), \alpha H(A^c))$$

For any  $A \in A$ , we have

$$E[P(A)] = \frac{\alpha H(A)}{\alpha H(A) + \alpha H(A^c)} = H(A)$$

$$\text{Var}[P(A)] = \frac{\alpha^2 H(A) H(A^c)}{(\alpha H(A) + \alpha H(A^c))^2 (\alpha H(A) + \alpha H(A^c) + 1)} = \frac{H(A)(1-H(A))}{1+\alpha}$$

Note:

$$P(P(A) - H(A) > \epsilon) \leq \frac{\text{Var}[P(A)]}{\epsilon^2} = \frac{H(A)(1-H(A))}{\epsilon^2(1+\alpha)}$$

Therefore,

$$P(A) \xrightarrow{d} H(A)$$

Imposing a large  $\alpha$  means that you are working more or less with  $H$ . Small  $\alpha$  means that, we provide more flexibility to  $P$ .

## Posterior of Dirichlet Process

If  $X_1, \dots, X_m$  is a set of realizations from  $P \sim DP(\alpha, H)$ , then

$$P(X_1, \dots, X_m \sim DP(\alpha_m, H_m))$$

where

$$\alpha_m = \alpha + m \quad \text{and} \quad H_m = \frac{\alpha}{\alpha+m} H + \frac{m}{\alpha+m} \sum_{i=1}^m \delta_{X_i}$$

$$\delta_X(A) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}$$

Therefore,  $H_m$  is a mixture of  $H$  and empirical process.

## Nonparametric Bayes Estimates

Under the squared error loss, an estimate for  $P$  is

$$E[P|X_1, \dots, X_m] = H_m^*$$

As  $m \rightarrow \infty$ ,  $H_m^* \xrightarrow{a.s.} P$ , where  $P$  is the actual distribution of  $X$ .

We also have

$$\sqrt{m} (P_m^* - P) \xrightarrow{d} B(F)$$

where  $B$  is Brownian bridge.

## Special case of Frullani Integral

$$\int_0^\infty (\exp(\theta u) - 1) \frac{\exp(-u)}{u} du = -\ln(1-\theta)$$

## Lévy Measure of Gamma( $\alpha, 1$ )

For  $\alpha > 0$ , let

$$N(x) = \alpha \int_x^\infty \frac{\exp(-t)}{t} dt$$

## Gamma Process

Define

$$G_t = \sum_{i=1}^\infty N^i(T_i + t)$$

$$T_i = E_1 + \dots + E_i \quad (E_i)_{i \geq 1} \stackrel{iid}{\sim} \text{Exp}(1)$$

The moment generating function for  $G_t$  is

$$M(\theta, t) = \exp\left(-\int_t^\infty (1 - \exp(\theta N^i(v))) dv\right)$$

$$M(\theta, 0) = (1-\theta)^{-\alpha}$$

Therefore, we get

$$\sum_{i=1}^\infty N^i(T_i) \sim T(\alpha, 1)$$

Let  $\Theta_i \stackrel{iid}{\sim} P_0$ . Then the set indexed process

$$G(\cdot) = \sum_{i=1}^\infty N^i(T_i) \delta_{\Theta_i}(\cdot)$$

is called a Gamma process

Theorem: We have

$$G(A) \sim T(\alpha P_0(A), 1)$$

$G(A)$  is independent of  $G(B)$  if  $A \cap B = \emptyset$

## Dirichlet Process (Formal Definition)

Let  $\Theta_i \stackrel{iid}{\sim} P_0$  be independent from  $\{T_i\}$  such that  $\alpha > 0$ . The random probability measure

$$P(\cdot) = \sum_{i=1}^\infty \frac{N^i(T_i)}{\sum_{i=1}^\infty N^i(T_i)} \delta_{\Theta_i}(\cdot)$$

is called a Dirichlet process with concentration parameter  $\alpha$  and concentration measure  $P_0$  and is denoted  $P \sim \text{Dir}(\alpha, P_0)$

Most of the time  $P_0$  is defined on  $\mathbb{R}^d$ . In this case:

$$H(x_1, \dots, x_n) = P_0(X_1 \leq x_1, \dots, X_n \leq x_n)$$

We usually write

$$P \sim \text{Dir}(\alpha, H)$$

The Dirichlet Process is difficult to use in practice since there's no closed form for  $N^i(x)$ . There are also infinitely many weights to it.

The above representation of the Dirichlet process is known as the Ferguson's Series Representation, where

$$\cdot T_i = E_1 + \dots + E_i, (E_i)_{i \geq 1} \stackrel{iid}{\sim} \text{Exp}(1)$$

$$\cdot (\Theta_i)_{i \geq 1} \stackrel{iid}{\sim} H$$

$\cdot (\Theta_i)_{i \geq 1}$  and  $(T_i)_{i \geq 1}$  are independent

$\cdot N(x)$  is the Lévy measure of a Gamma( $\alpha, 1$ ) random variable.

## Sethuraman's Representation of DP



$$P_{\text{seth}}(\cdot) = \sum_{i=1}^{\infty} P_i \delta_{x_i}(\cdot)$$

where:

- $(\theta_i)_{i \geq 1}$  iid H
  - $(P_i)_{i \geq 1}$  is an independent sequence of random variables (called weights) such that
- $$P_i = \theta_i, P_i = \frac{1}{\prod_{k=1}^i (1-\theta_k)}, i \geq 2$$
- $(P_i)_{i \geq 1}$  iid Beta( $1, \alpha$ )
  - $(\theta_i)_{i \geq 1}$  and  $(P_i)_{i \geq 1}$  independent

## Sethuraman's Truncated Representation

$$P_n^{\text{seth}}(\cdot) = \sum_{i=1}^n P_i \delta_{x_i}(\cdot)$$

where:

- $(\theta_i)_{i \geq n}$  and  $(P_i)_{i \geq n}$  are defined as in Sethuraman's representation with  $P_n=1$
  - The assumption that  $P_n=1$  is necessary to make the weights add to 1 almost surely
  - A random stopping rule for choosing  $n=n(\epsilon)$  was proposed by Muliere and Tradella (1998), where, for each  $\epsilon$ ,
- $$n=\inf\{i: P_i=(1-\beta_1)\dots(1-\beta_{i-1}), \beta_i < \epsilon\}$$

## Zarepour's Theorem

Let  $(P_{1,n}, \dots, P_{n,n}) \sim \text{Dir}(\frac{1}{n}, \dots, \frac{1}{n})$  and  $X_1, X_2, \dots \stackrel{iid}{\sim} H$ . Then

$$\sum_{i=1}^n P_{i,n} \delta_{X_i} \xrightarrow{d} DP(\alpha, H)$$

## Irreducible Markov Chains

A Markov Chain is irreducible if  $i \leftrightarrow j$ , for all  $i, j \in S$ . In other words, there is only one class S.

## Aperiodic Markov Chains

We first need a definition: Let  $X_n$  be a Markov Chain with state space  $\{1, 2, \dots, n\}$ . Let  $i \in S$ ,  $d(i) = \text{g.c.d.}(n, P_{ii}^{(n)})$

A Markov Chain  $X_n$  is called aperiodic if  $d(i)=1$  for all  $i \in S$ .

We can also show that if  $i \leftrightarrow j$  then  $d(i)=d(j)$ .

## Recurrent and Transient States

We say a state  $i$  is recurrent if  $\sum_{n=0}^{\infty} P_{ii}^{(n)} = \infty$ . Otherwise, if  $\sum_{n=0}^{\infty} P_{ii}^{(n)} < \infty$ , the state is called transient.

Define

$$f_{ij}^{(n)} = P(X_n=i, X_j \neq i, j \neq 1, \dots, n-1 | X_0=i)$$

Therefore,

$$P_{ii}^{(n)} = \sum_{j \neq i} f_{ij}^{(n)} P_{ji}^{(n)}$$

$$f_{ii}^{(n)} = 0, f_{ii}^{(n)} = P_{ii}^{(n)}, P_{ii}^{(n)} = 1$$

If we define the generating functions for  $\{f_{ij}^{(n)}\}$  and  $\{f_{ii}^{(n)}\}$  by

$$P_{ij}(s) = \sum_{n=0}^{\infty} P_{ij}^{(n)} s^n$$

$$F_{ii}(s) = \sum_{n=0}^{\infty} f_{ii}^{(n)} s^n$$

where  $f_{ij}^{(n)}$  is the probability that the first passage from state  $i$  to state  $j$  occurs at the  $k$ -th transition. We can see that

$$F_{ii}(s) P_{ij}(s) = P_{ij}^{(n)} - 1$$

Therefore,

$$P_{ii}(s) = \frac{1}{1 - F_{ii}(s)}$$

The state  $i$  is transient if  $\sum_{n=0}^{\infty} P_{ii}^{(n)} < \infty$ . Therefore  $\lim_{n \rightarrow \infty} P_{ii}^{(n)} = 0$ .

The state  $i$  is null recurrent if  $\sum_{n=0}^{\infty} P_{ii}^{(n)} = \infty$  but  $\lim_{n \rightarrow \infty} P_{ii}^{(n)} = 0$ .

The state  $i$  is positive recurrent if  $\sum_{n=0}^{\infty} P_{ii}^{(n)} = \infty$  and  $\lim_{n \rightarrow \infty} P_{ii}^{(n)} = \frac{1}{\pi_i}$ .

A Markov chain is ergodic if it is irreducible, aperiodic, and all states are positive recurrent.

## Time Reversible Markov Chains

When a Markov chain is at a steady state, the total probability transferring out of state  $j$  must balance the total amount of probability transferring into state  $j$ . This must be true for all states  $j$ .

We can view the states in reverse time. If  $Q$  is the reverse probability transition matrix, we have

$$Q_{ij} = \frac{P(X_{n+1}=j | X_n=i)}{P(X_{n+1}=i)}$$

## At the steady state this becomes

$$Q_{ij} = \frac{\pi_j P_{ji}}{\pi_i}$$

The Markov chain is reversible if the backward Markov chain and the forward Markov chain have the same transition probabilities. In other words,  $P=Q$ .

In this case, for all  $i$  and  $j$

$$\pi_i P_{ij} = \pi_j P_{ji}$$

This is called detailed balance.

Note that  $\pi_i P_{ij}$  is the probability going from state  $i$  to state  $j$  and  $\pi_j P_{ji}$  is the probability going from state  $j$  to state  $i$ . The detailed balance says that, at the steady state for any two states  $i$  and  $j$  the transfer between states balance each other.

Since  $\pi_i P_{ij} = \pi_j P_{ji}$  for all  $i$ , the transition from any state to itself always satisfies the detailed balance.

Theorem: If the probability matrix  $P$  for a Markov chain satisfies detailed balance condition then we have

$$\pi P = \pi$$

## Metropolis Algorithm

We will use accept reject method to accept some transitions to get desired steady state. Notice that if detailed balance for states  $i$  and  $j$  does not satisfy then

$$(i) \pi_i P_{ij} < \pi_j P_{ji} \quad \text{or} \quad (ii) \pi_i P_{ij} > \pi_j P_{ji}$$

If (i) holds, we accept all transitions from state  $i$  to state  $j$  and only accept all the transitions from  $j$  to  $i$  to get the steady balance. Similarly when (ii) holds, then we accept only some of the transitions from state  $i$  to  $j$ . We need extra probability to the transition from  $i$  to  $j$  to compensate for the unaccepted transitions.

To start the algorithm, we start from the transition probabilities with the desired steady state probabilities  $\pi_i$ .

Then

• Define the acceptance probability

$$\alpha_{ij} = \min\left(\frac{\pi_j P_{ji}}{\pi_i P_{ij}}, 1\right), \forall i, j$$

• For each  $i$  and  $j \neq i$ , let  $P_{ij}' = \alpha_{ij} P_{ij}$

• Let  $P_{ii}' = 1 - \sum_{j \neq i} P_{ij}'$

Then  $\pi' = (\pi_1, \pi_2, \dots)$  is the steady state distribution for the Markov chain with transition matrix  $P' = (P'_{ij})$

Theorem: The Markov chain with transition probability matrix  $P'$  defined above satisfies the detailed balanced condition and has the desired steady state distribution  $\pi'$ .

## Continuous State Markov Chains

In a Markov process with continuous distribution we define the transition model by

$$P(x, A) = P(X_{n+1} \in A | X_n=x)$$

In this case, the transition c.d.f. is

$$F(v|x) = P(X_{n+1} \leq v | X_n=x)$$

and its transition p.d.f. is

$$f(v|x) = \frac{dF(v|x)}{dx}$$

We can write in general

$$f^{(n)}(v|x) = \int_{-\infty}^v f^{(n)}(w|x) f^{(n)}(w|x) dw$$

If  $\lim_{n \rightarrow \infty} f^{(n)}(v|x) = g(v)$ , we can write

$$g(v) = \int_{-\infty}^v g(w)f(v|w) dw$$

A time reversible Markov chain satisfies the detailed balance condition. If  $g$  is the steady state p.d.f. then we can write

$$g(x)f(v|x) = g(v)f(x|v)$$

**T H E  
E N D**

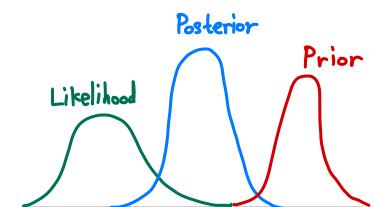
You Made it

MAT 4381

Winter 2024

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$h(\theta|y) = \frac{f(y|\theta) g(\theta)}{\int f(y|\theta) g(\theta) d\theta}$$



Good Luck on the exam!!

You got this!!!