

# Statistical Inference

## Lecture 10b

ANU - RSFAS

Last Updated: Tue May 8 16:45:05 2018

# Bayesian Testing

- Consider testing:

$$H_0 : \theta \in \Theta_0 \quad \text{vs.} \quad H_1 : \theta \in \Theta_1$$

- The classical statistician considers  $\theta$  a **fixed** unknown, thus a hypothesis test is **true** or **false**. Either  $\theta$  is in  $\Theta_0$  or  $\Theta_1$ !
- Bayesians however consider  $\theta$  to be **random** and it is quite natural to consider:

$$P(\theta \in \Theta_0 | \mathbf{x}) \quad \text{vs.} \quad P(\theta \in \Theta_1 | \mathbf{x})$$

# Bayesian Testing

- One approach to Bayesian testing is to reject  $H_0$  if:

$$\underline{P(\theta \in \Theta_1 | \mathbf{x}) > P(\theta \in \Theta_0 | \mathbf{x})}$$

# Bayesian Testing

**Example 8.2.7:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{normal}(\theta, \sigma^2)$ . Let  $\theta \sim \text{normal}(\mu, \tau^2)$ , where  $\sigma^2, \mu, \tau^2$  are known. Consider testing:

$$H_0 : \theta \leq \theta_0 \quad \text{vs.} \quad H_1 : \theta > \theta_0$$

$$[\theta|\mathbf{x}] \sim \text{normal} \left( \frac{\sigma^2 \mu + n\tau^2 \bar{x}}{\sigma^2 + n\tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} \right) \quad \text{normal posterior}$$

- To compare we simply examine:

$$\int_{-\infty}^{\theta_0} p(\theta|\mathbf{x}) d\theta \quad \text{vs.} \quad \int_{\theta_0}^{\infty} p(\theta|\mathbf{x}) d\theta$$

# Bayesian Testing

- Based on this Bayesian approach to hypothesis testing, what if we wanted to test:

$$H_0 : \theta = \theta_0 \quad \text{vs.} \quad H_1 : \theta \neq \theta_0$$

$$[\theta|\mathbf{x}] \sim \text{normal} \left( \frac{\sigma^2 \mu + n\tau^2 \bar{x}}{\sigma^2 + n\tau^2}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} \right)$$

- $[\theta|\mathbf{x}]$  is a continuous distribution, so the probability of any single point (such as  $\theta_0$ ) is zero. This approach does not seem to work.

# Bayesian Testing



- Notice that when we test through a Bayesian approach, we model the parameter of interest (i.e. we put a prior on it).
- If our scientific question concerns whether a parameter can be **exactly  $\theta_0$  or not**, then we should model that:

*"prior is a mixture"*

$$\theta \sim p\mathbf{1}_{\theta=\theta_0} + (1-p)\text{normal}(\mu, \tau^2)$$

where  $0 \leq p \leq 1$ .

- In a regression or GLM setting many times we are interested in testing whether  $\beta = 0$  vs.  $\beta \neq 0$  and these priors or variants on them can be quite useful.

# Bayesian Testing

②

- Another approach that does allow for consideration of  $\beta = 0$  is Bayes factors. Let's rephrase hypothesis testing as choosing between competing models:

$$\left\{ \begin{array}{ll} \text{Model 1 } (M_1): y_i = \alpha + \epsilon_i & \epsilon_i \stackrel{\text{iid}}{\sim} \text{normal}(0, \sigma^2); \quad \theta_1 = \{\alpha, \sigma^2\} \\ \text{Model 2 } (M_2): y_i = \alpha + \beta x_i + \epsilon_i & \epsilon_i \stackrel{\text{iid}}{\sim} \text{normal}(0, \sigma^2); \quad \theta_2 = \{\alpha, \beta, \sigma^2\} \end{array} \right.$$

# Bayesian Testing

*based on p.p. choose a model.*

- Let's figure out the posterior probability for a given model  $i$ :

*marginal prob. of the data given the model.*

$$\pi(M_i|\mathbf{x}) = \frac{f(\mathbf{x}|M_i)\pi(M_i)}{m(\mathbf{x})}$$

*→ marginal prob of the data.*

$$f(\mathbf{x}|M_i) = \int_{\theta} f(\mathbf{x}|\theta, M_i)\pi(\theta|M_i)d\theta_i$$

$$m(\mathbf{x}) = \sum_{i=1}^2 f(\mathbf{x}|M_i)\pi(M_i)$$



# Bayesian Testing

- Now consider the following ratio of the posterior model probabilities:

$$\begin{aligned}\frac{\pi(M_2|\mathbf{x})}{\pi(M_1|\mathbf{x})} &= \frac{f(\mathbf{x}|M_2)}{f(\mathbf{x}|M_1)} \times \frac{\pi(M_2)}{\pi(M_1)} \\ &= BF(M_2; M_1) \times \frac{\pi(M_2)}{\pi(M_1)}\end{aligned}$$

*interest!*

Where  $BF(M_2; M_1)$  is called the Bayes factor.

- Typically  $\pi(M_2) = \pi(M_1)$ , so the ratio of the posterior probabilities is the Bayes factor.
- The Bayes factor looks like a likelihood ratio. However, the difference is that  $\theta$  has been integrated out in both the numerator and denominator, so we have the marginal distribution of the data given the model.

# Bayesian Testing

- If  $f(\mathbf{x}|M_2) > f(\mathbf{x}|M_1)$  or  $\frac{f(\mathbf{x}|M_2)}{f(\mathbf{x}|M_1)} > 1$  then we have support for  $M_2$  against  $M_1$ .
- Jeffreys, H. (1961 appendix B) suggested the following:

$BF(M_2; M_1) = B_{21}$	Evidence against model 1 ( $H_0$ )
1 to 3.2	Not worth more than a bare mention
3.2 to 10	Substantial
10 to 100	Strong
> 100	Decisive

# Bayesian Interval Estimation

- Suppose we have data  $X_1, \dots, X_n$  from density  $f_X(x|\theta)$  along with a prior distribution  $\pi(\theta)$ . As we saw we use Bayes' rule to update our 'beliefs' about  $\theta$  once we observe the data:

$$\begin{aligned}\pi(\theta|\mathbf{x}) &= \frac{L(\theta|\mathbf{x})\pi(\theta)}{\int_{\theta \in \Theta} L(\theta|\mathbf{x})\pi(\theta)d\theta} \\ &= \frac{L(\theta|\mathbf{x})\pi(\theta)}{m(\mathbf{x})}\end{aligned}$$

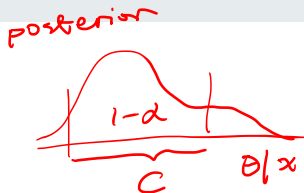
- So we have the whole distribution for

$$\pi(\theta|\mathbf{x})$$

- This is different than the frequentist approach where find an estimator for  $\theta$ , say  $\hat{\theta}$  and then try to determine the distribution of  $\hat{\theta}$ .

# Bayesian Interval Estimation

- To obtain an interval we simply consider:



$$P\pi(\theta|x)(C) = \int_C \pi(\theta|x) d\theta = 1 - \alpha$$

- Be careful. We are using  $\alpha$  quite generically. Recall that  $\alpha$  does have a formal definition: The probability of a Type-I error. This is based on repeated sampling. For the Bayesian case we only think about one data set an infinite number of possible data sets.
- There are quite a lot of choices for  $C$ . We will consider the 3 most common.

# Bayesian Interval Estimation

- ①
- Equal tailed: both sides with  $\frac{\alpha}{2}$  tails

$$\int_{-\infty}^{\theta_L} \pi(\theta|\mathbf{x})d\theta = \alpha/2, \quad \int_{\theta_U}^{\infty} \pi(\theta|\mathbf{x})d\theta = \alpha/2$$

- ②
- Smallest length: We can choose  $C$  to minimize  $\theta_U - \theta_L$ .

# Bayesian Interval Estimation

③

- Highest posterior density region (HPD): We define  $C$  to be that set with posterior probability  $1 - \alpha$  which satisfies the criterion:

$$\theta_1 \in C \quad \text{and} \quad \overset{\text{"prob"}}{\pi(\theta_2|\mathbf{x})} > \pi(\theta_1|\mathbf{x}) \Rightarrow \theta_2 \in C$$

$C$  contains the values of  $\theta$  which have the highest posterior density values, so that we can determine HPD regions as the set:

$$C = \{\theta \in \Theta : \pi(\theta|\mathbf{x}) > c_\alpha\}$$

- If the posterior is unimodal then this will be the smallest length interval!

"centred around the peak"

# Bayesian Interval Estimation

$$E(X) = \frac{1}{\theta}$$

**Example:**  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{exponential}(1/\theta)$  and  $\pi(\theta) = \theta \exp(-\theta)$ .

$$\begin{aligned}\pi(\theta|\mathbf{x}) &\propto \underbrace{\left\{ \prod_{i=1}^n \theta \exp(-x_i \theta) \right\}}_{\text{likelihood}} \underbrace{\theta \exp(-\theta)}_{\text{prior}} \\ &= \theta^n \exp\left(-\sum x_i \theta\right) \theta \exp(-\theta) \\ &= \theta^{n+1} \exp(-\theta(n\bar{x} + 1)) \\ &= \theta^{n+2-1} \exp(-\theta(n\bar{x} + 1))\end{aligned}$$

posterior  
is a Gamma

$$[\theta|\mathbf{x}] \sim \text{gamma}\left(n+2, \frac{1}{n\bar{x}+1}\right)$$

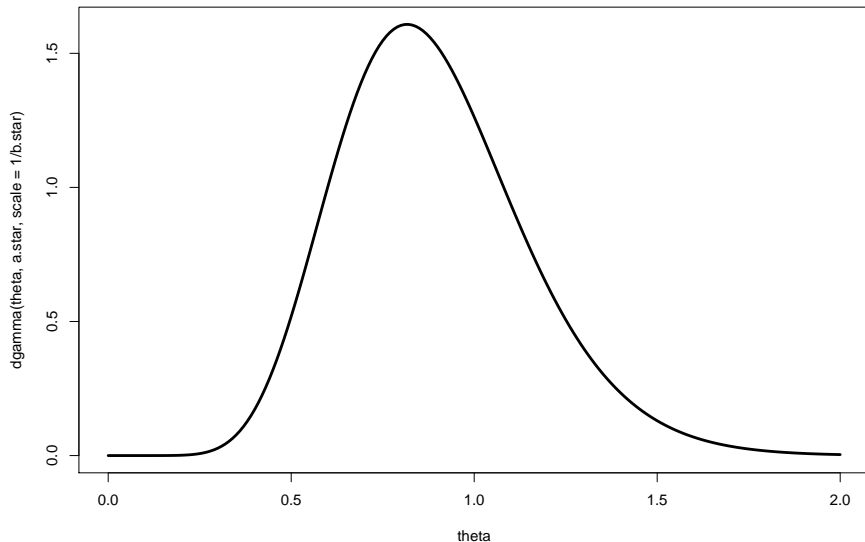
# Bayesian Interval Estimation

- Let's plot the density for  $n = 10$  and  $\bar{x} = 1.247$ .

```
n=10; x.bar <- 1.247
a.star <- n+2; b.star <- n*x.bar+1
theta <- seq(0, 2, by=0.01)
plot(theta, dgamma(theta, a.star,
                    scale=1/b.star), type="l", lwd=3)
```



# Bayesian Interval Estimation



# Bayesian Interval Estimation

- An equal-tailed 95% interval is given by  $[\theta_l, \theta_u]$ :

$$\begin{aligned}\int_0^{\theta_l} \pi(\theta|\mathbf{x}) &= 0.025 \\ F_{[\theta|\mathbf{x}]}(\theta_l) &= 0.025\end{aligned}$$

$$\begin{aligned}\int_0^{\theta_u} \pi(\theta|\mathbf{x}) &= 1 - 0.025 = 0.975 \\ F_{[\theta|\mathbf{x}]}(\theta_u) &= 0.975\end{aligned}$$

# Bayesian Interval Estimation

```
theta.L <- qgamma(0.025, a.star, scale=1/b.star)
theta.U <- qgamma(0.975, a.star, scale=1/b.star)

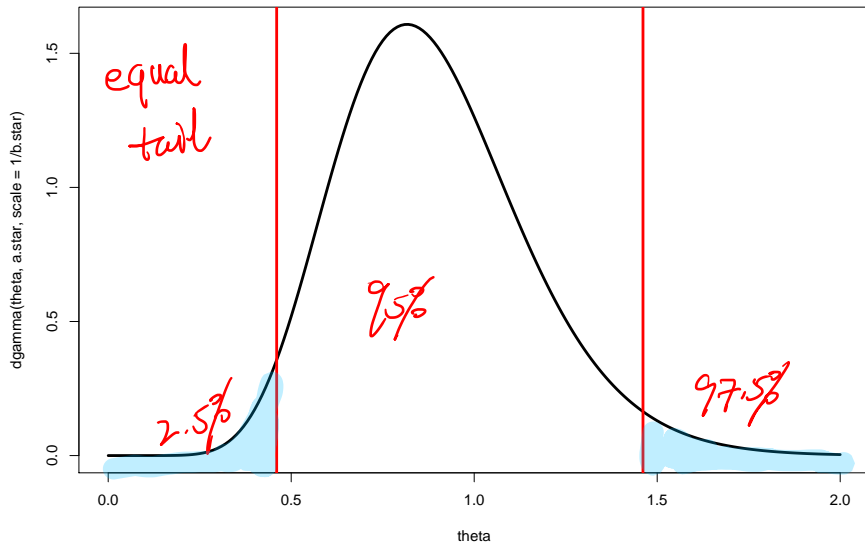
c(theta.L, theta.U)
```

```
## [1] 0.4603248 1.4611758
```

*"credible interval"*

```
plot(theta, dgamma(theta, a.star, scale=1/b.star), type="l", lwd=2, col="blue",
      abline(v=c(theta.L, theta.U), lwd=3, col="red"))
```

# Bayesian Interval Estimation



# Bayesian Interval Estimation

- If we only have tables in front of us, we can relate the gamma distribution to a  $\chi^2$  distribution as was discussed in tutorial the other week.
- If  $[\theta|\mathbf{x}] \sim \text{gamma}(a^*, b^*)$  then

$$\left[ \frac{2\theta}{b^*} \middle| \mathbf{x} \right] \sim \text{gamma}(a^*, 2)$$

$$\sim \chi^2_{p=2a^*}$$

# Bayesian Interval Estimation

- Using probabilities to the left.  $p = 2a^* = 2n + 4$ .

pivoting

$$\begin{aligned} \left[ \chi_{0.025,p}^2 \leq \frac{2\theta}{b^*} \middle| \mathbf{x} \leq \chi_{0.975,p}^2 \right] \\ \left[ \chi_{0.025,p}^2 \leq 2\theta(n\bar{x} + 1) \middle| \mathbf{x} \leq \chi_{0.975,p}^2 \right] \\ \left[ \frac{\chi_{0.025,p}^2}{2(n\bar{x} + 1)} \leq \theta \middle| \mathbf{x} \leq \frac{\chi_{0.975,p}^2}{2(n\bar{x} + 1)} \right] \end{aligned}$$

```
p <- 2*n + 4
theta.L <- qchisq(0.025, p)/(2*(n*x.bar+1))
theta.U <- qchisq(0.975, p)/(2*(n*x.bar+1))

c(theta.L, theta.U)
```

```
## [1] 0.4603248 1.4611758
```

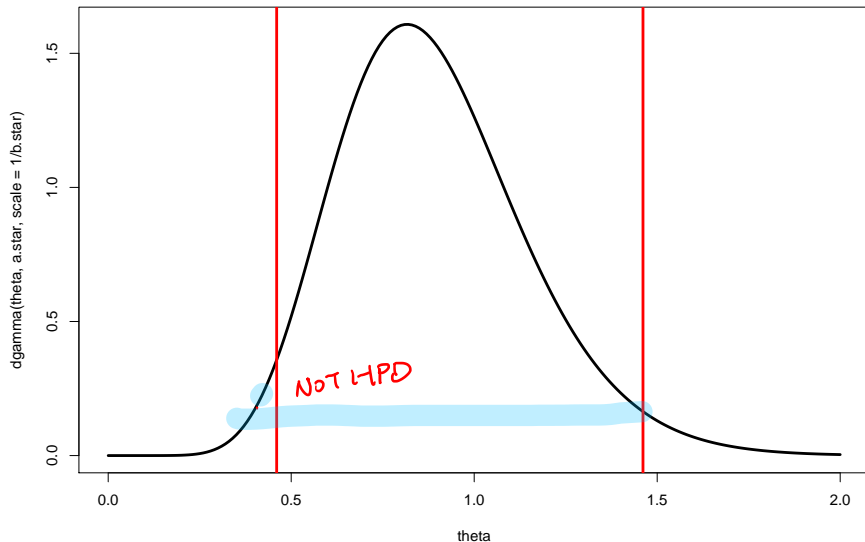
# Bayesian Interval Estimation

- Is the interval  $[0.4603, 1.4612]$  a HPD (highest posterior density) interval (the posterior is unimodal)?
- Recall:

$$\theta_1 \in C \quad \text{and} \quad \pi(\theta_2|\mathbf{x}) > \pi(\theta_1|\mathbf{x}) \Rightarrow \theta_2 \in C$$

- Let's see the density with the equal-tailed interval again.

# Bayesian Interval Estimation





# Bayesian Interval Estimation

- Note that the density seems to be higher for  $\theta = 0.40$  than  $\theta = 1.4612$ :

```
dgamma(0.4, a.star, scale=1/b.star)
```

```
## [1] 0.1713707
```

```
dgamma(1.4612, a.star, scale=1/b.star)
```

```
## [1] 0.1641042
```

- So the equal-tailed interval is not a HPD interval!

*density should be the same!*  
*leftmost & rightmost "tails" are*

# Bayesian Interval Estimation

- To get the HPD interval we take horizontal slices across the density till we get the appropriate probability.

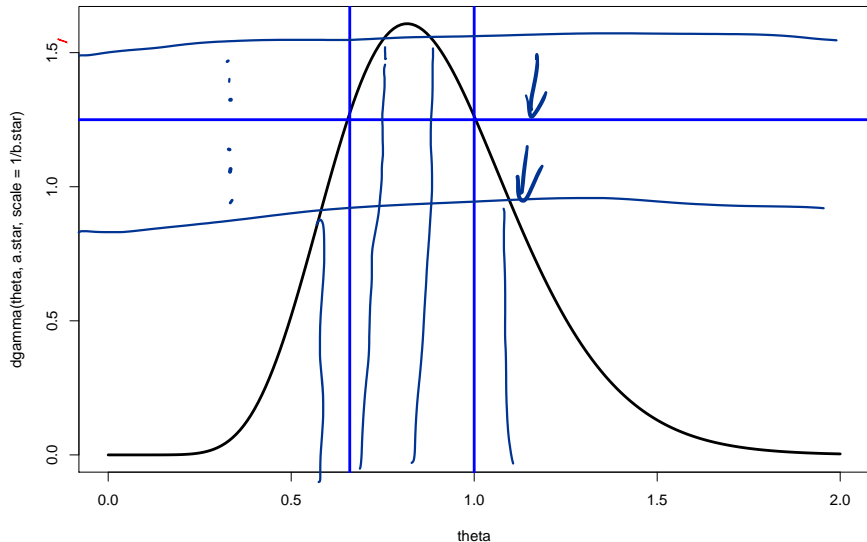
```
plot(theta, dgamma(theta, a.star, scale=1/b.star), type="l", lwd=3)
abline(h=1.25, lwd=3, col="blue")
```

```
##
theta <- seq(0, 2, by=0.01)
dens <- dgamma(theta, a.star, scale=1/b.star)
```

```
##
hpd.cut <- 1.25
theta.L <- min(theta[dens>=hpd.cut])
theta.U <- max(theta[dens>=hpd.cut])
abline(v=c(theta.L, theta.U), lwd=3, col="blue")
```

```
## interval probability
pgamma(theta.U, a.star, scale=1/b.star) -
  pgamma(theta.L, a.star, scale=1/b.star)
```

# Bayesian Interval Estimation



```
## [1] 0.5062717
```

# Bayesian Interval Estimation

```
hpd.cut <- sort(seq(0.1, 1.25, by=0.0001), decreasing =TRUE)
c <- 1
cred.int <- 0.5063

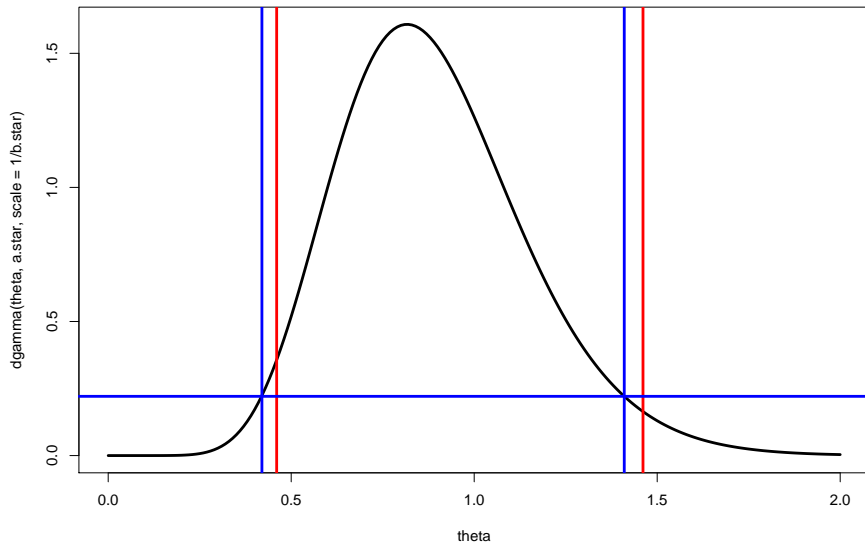
while(cred.int<0.95){
  theta.L <- min(theta[dens>=hpd.cut[c]])
  theta.U <- max(theta[dens>=hpd.cut[c]])

  ## interval probability
  cred.int <- pgamma(theta.U, a.star, scale=1/b.star) -
    pgamma(theta.L, a.star, scale=1/b.star)
  c <- c+1
}

HPD <- c(theta.L,theta.U)
HPD

## [1] 0.42 1.41
```

# Bayesian Interval Estimation



# Bayesian Interval Estimation

- Let's check the length of each interval:
  - equal-tailed:  $1.46 - 0.460 = 1.00$
  - HPD:  $1.41 - 0.42 = 0.99$
- HPD is the shorter interval, but not by much.

# Bayesian Inference: Properties

**Definition 7.1:** A statistic  $T(\mathbf{X})$  is sufficient for  $\theta$  if and only if the posterior distribution of  $\theta$  given  $\mathbf{X}$  is the same as the posterior distribution of  $\theta$  given  $T(\mathbf{X})$ .

**Proof:** Note that Definitions 2.5 and 7.1 are the same!

Suppose that  $T(\mathbf{X})$  satisfies Definition 2.5. Then:

no longer have theta there

$$f(\mathbf{x}; \theta) = g(\mathbf{x}|t, \theta)h(t|\theta) = g(\mathbf{x}|t)h(t|\theta)$$

- The posterior is

$$\begin{aligned} p(\theta|\mathbf{x}) &\propto f(\mathbf{x}; \theta)p(\theta) \\ &\propto h(t|\theta)p(\theta) \\ &\propto p(\theta|t) \end{aligned}$$

- Now assume that  $T(\mathbf{X})$  satisfies Definition 7.1.

$$\begin{aligned} f(\mathbf{x}|\theta) &= \frac{p(\theta|\mathbf{x})h(\mathbf{x})}{p(\theta)} \\ \text{likelihood} &= \frac{p(\theta|\mathbf{x})h(\mathbf{x})}{p(\theta)} \\ &= K_1[\mathbf{x}|\theta] K_2[\mathbf{x}] \end{aligned}$$

- From the **factorization theorem**, it follows that  $T(\mathbf{X})$  is a sufficient statistic.



# Bayesian Inference: Asymptotics - Rough Idea

- Suppose we have  $y_1, \dots, y_n \sim p(y|\theta)$ .
- Let's consider the posterior distribution:

$$\begin{aligned} p(\theta|\mathbf{y}) &\propto p(\mathbf{y}|\theta) p(\theta) \\ &= \exp[\log p(\mathbf{y}|\theta)] \exp[\log p(\theta)] \end{aligned}$$

- As  $n \rightarrow \infty$  the posterior is dominated by the likelihood.

$$p(\theta|\mathbf{y}) \propto \exp[\log p(\mathbf{y}|\theta)]$$

- Thus to an approximation we have the following:

$$\begin{aligned}
 p(\theta|\mathbf{y}) &\propto \exp[\log p(\mathbf{y}|\theta)] \\
 &\propto \exp[\ell(\theta)] \\
 &\propto \exp\left[\ell(\hat{\theta}) + \underbrace{(\theta - \hat{\theta})}_{\text{red } 0} \underbrace{\ell'(\hat{\theta})}_{\text{red } 0} + \frac{1}{2}(\theta - \hat{\theta})^2 \ell''(\hat{\theta})\right] \\
 &\propto \exp\left[\frac{1}{2}(\theta - \hat{\theta})^2 \ell''(\hat{\theta})\right]
 \end{aligned}$$

- Where:  $\hat{\theta}$  is the MLE.
- Note:  $\ell(\hat{\theta})$  is a constant.
- Note:  $\ell'(\hat{\theta}) = 0$

$X \sim P(\theta)$   
 $E[X] \rightarrow \theta$   
 $\frac{1}{N} \sum_{i=1}^N X_i \rightarrow E[\theta] = \theta$   
 sample mean  
 dist mean

$$\begin{aligned}
 p(\theta|\mathbf{y}) &\propto \exp \left[ \frac{1}{2} (\theta - \hat{\theta})^2 \ell''(\hat{\theta}) \right] \quad \text{converges to that} \\
 &\propto \exp \left[ \frac{1}{2} (\theta - \hat{\theta})^2 [-I(\hat{\theta})] \right] \quad I''(\hat{\theta}) = H(\hat{\theta}) \\
 &\propto \exp \left[ - \frac{1}{2 [I(\hat{\theta})]^{-1}} (\theta - \hat{\theta})^2 \right] \\
 &\quad \text{the kernel for normal distribution}
 \end{aligned}$$

- We see that this expression is proportional to a normal distribution. So we have:

$$p(\theta|\mathbf{y}) \approx \text{normal} \left( \hat{\theta}, [I(\hat{\theta})]^{-1} \right)$$