

# **BAYESIAN INFERENCE**

**BY**

**TANUJIT CHAKRABORTY**

**Indian Statistical Institute**

**Mail : [tanujitisi@gmail.com](mailto:tanujitisi@gmail.com)**

# BAYESIAN INFERENCE

$f(x|\theta)$   $x$ : Effect  $\theta$ : Cause : Probability Thinking  
 $l(\theta|x)$  Inversion process.

Bayesian thinking  $\leftrightarrow$  Inverse probability thinking.

Historical Background: - Probability has been a subject of study for a long time. Statistics is relatively a young field. Linear regression due to Galton appeared in the late 1800's and Karl Pearson's goodness of fit measure and correlation appeared in 19th century.

The field of statistics blossomed in the 1920's and 1930's with works of Fisher, Neyman, Pearson. A flurry of research was prompted by World War II, which generated a wide variety of difficult applied problems and the first substantive government funding came in US and UK.

By contrast, Bayesian methods are much older, dating to the original paper of 'THOMAS BAYES' in the 1760s. The area generated some interest in the works of Laplace, Gauss and others in the 19th century. The Bayesian approach was ignored and opposed by the statisticians of the early 20th century. Fortunately, during this period, several prominent non-statisticians, most notably Harold Jeffreys (a Physicist) and Arthur Bowley (an Econometrician) continued to lobby on behalf of Bayesian ideas (which they referred to as "Inverse Probability").

Beginning 1950's, statisticians such as L.J. Savage, Bruno de Finetti, Dennis Lindley and many others come in and they pointed out several deficiencies of the classical approach. The biggest impetus to Bayesian statistics came in 1990's after many computational algorithms and modern computers started to surface.

Motivating Examples: - "Acceptance of paper in Anals of Statistics".

What is the probability of accepting your paper with a priori information that your 5 papers submitted to Anals got accepted?

Some Paradoxes:

• Example 1:  $X_1, X_2 \stackrel{iid}{\sim} U\left[\theta - \frac{1}{2}, \theta + \frac{1}{2}\right], \theta \in \mathbb{R}$ .  
A 95% CI for  $\theta$  is given by  $[\bar{x} - c, \bar{x} + c]$  by choosing  $c$  properly.

Now let us suppose we are given  $X_{(2)} = 2$  and  $X_{(1)} = 1$ .

Then  $\theta$  is obviously  $\frac{3}{2} = \bar{x}$ .

But a frequentist, if he sticks to the long term frequency interpretation of CI, will not be able to say if the given interval contains  $\theta$  or not surely !!

• Example 2: (Cox 1958)

To estimate  $\mu$  in  $N(\mu, \sigma^2)$ , toss a fair coin. Draw a random sample of size  $n=2$  if it falls heads and draw a sample of size  $n=1000$  if it falls tails. An unbiased estimate of  $\mu$  is  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  with variance  $\frac{1}{2} \left( \frac{\sigma^2}{2} + \frac{\sigma^2}{1000} \right) \approx \frac{\sigma^2}{4}$ .

Will you believe that the variance is like  $\frac{\sigma^2}{4}$  if you actually draw a sample of size 1000 and compute  $\bar{X}$  based on that sample?

• Example 3: (Lindly and Philips (1976)).

Suppose in 12 independent tosses of a coin you observe 9 heads and 3 tails and the last toss producing a tail.

To test  $H_0: \theta = \frac{1}{2}$  Vs.  $H_1: \theta > \frac{1}{2}$  ( $\theta$ : Prob. of success)

We apply two choices of possible sampling distributions.

(a) Binomial: Here  $n=12$  tosses are fixed before hand and the random quantity  $X$  is the # of heads obtained in  $n$  tosses.

Then  $X \sim \text{Bin}(12, \theta)$ .

Likelihood (for the given data) is given by  $\binom{12}{9} \theta^9 (1-\theta)^3$ .

(b) Negative Binomial: Here one has to assume that data collection continued until this tail appeared. The random quantity  $X$  is the number of heads required to complete the experiment so that  $X \sim \text{Neg. Bin}(r=3, \theta)$ .

Likelihood is given by  $\binom{11}{9} \theta^9 (1-\theta)^3$ .

Consider the Rejection Region:  $\text{Reject } H_0 \text{ if } X \geq c.$

p-value:-

Binomial case:  $\alpha_1 = P_{\theta = \frac{1}{2}}(X \geq 9) = \sum_{j=9}^{12} \binom{12}{j} \theta^j (1-\theta)^{12-j} = 0.075$

Neg-Bino. case:  $\alpha_1 = P_{\theta = \frac{1}{2}}(X \geq 9) = \sum_{j=9}^{\infty} \binom{12+j}{j} \theta^j (1-\theta)^3 = 0.0325.$

Using the "usual" Type I error level  $\alpha = 0.05$ , two models lead to two decisions. You should note that the likelihoods are propatrical here.

One of the main thing what went 'wrong' is that we allowed unobserved outcomes to affect the rejection decision. That is, probability of  $X$  values may be greater than 9 (the value actually observed) was used as evidence against  $H_0$  in each case, even though these values did not occur.

As mentioned before, statistical inference is based on a probabilistic modelling of the observed phenomenon. There are two approaches:

- A first approach assumes that statistical inference must incorporate as much as possible of the complexity of the phenomenon, and thus aims at estimating the distributions underlying the phenomenon under minimal assumptions, generally using functional estimation (density, regression, etc). This is non-parametric approach.
- conversely, the parametric approach represents the distribution through a density function  $f(x|\theta)$ , where only the parameter  $\theta$  (of finite dimension) is unknown. This is the parametric approach.

- Reference Books:- 1. An Introduction to Bayesian Analysis by JK Ghosh, Mohan Delampady and Tapas Samanta
- 2. Bayesian Data Analysis 3 by Rubin, Stern, Goltman, Dunson
- 3. Monte Carlo Statistical Methods by Robert Casella.
- 4. Collected works of D. Basu by JK Ghosh, revised version by Anirban Dasgupta.

## Classical Inference Model:-

Let us consider for example, the problem of estimation of a real parameter.

$\theta$  is unknown in nature.

Data:  $\underline{X} \sim P_\theta, \theta \in \Theta$

$T(\underline{X})$ : estimator

For a good estimator,  $T(\underline{X})$ ,  $|T(\underline{X}) - \theta|$  should be small on equivalently  $(T(\underline{X}) - \theta)^2$  should be small.

Classical statistics averages this quantity over the whole sample space  $\mathcal{X}$  (set of all possible values of  $\underline{X}$ )

$$R(T, \theta) = E_\theta (T(\underline{X}) - \theta)^2$$

A Bayesian won't do this. He/she will consider a data dependent measure of performance of an estimator — something that depends on the particular data in hand, the data he has observed. Once he has a particular data that he has observed, he is not going to average over all other possible values of the data that he have not observed.

Classical Statistics:- An estimator  $T_0$  is said to be the best if it minimizes  $R(T, \theta)$  w.r.t. choices of the estimator  $T$ . Here  $\theta$  is unknown. So, this is a very difficult problem. Indeed, there does not exist a "best" estimator.

There are two ways:-

(i) Imposing restrictions to the class of estimators (such as unbiasedness and invariance) and finding the best estimator in the restricted class.

(ii) Considering all the estimators but evaluating them in a different manner.  
e.g.  $\sup_{\theta} R(T, \theta)$  (leading to the minimax estimator)

Example:- Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1); \theta \in \mathbb{R}$

UMVUE is  $\bar{X}$ . Suppose we know  $-1 \leq \theta \leq 1$ . Then UMVUE is a bad estimator. There are natural common sense estimators of  $\theta$  better than  $\bar{X}$ .

The Bayesian Model:- The purpose of statistical analysis is fundamentally an inversion purpose, since it aims at retrieving the causes — reduced to the parameter of the probabilistic generating mechanism from the effects — summarized by the observations.

In other words, when observing a random phenomenon directed by a parameter  $\theta$ , statistical methods allow to deduce from these observations an inference (i.e., a summary of a characterization) about  $\theta$ .

A general description of inversion of probabilities is given by Bayes' formula.

If  $A$  and  $E$  are events such that  $P(E) \neq 0$ ,  $P(A|E)$  and  $P(E|A)$  are related by

$$P(A|E) = \frac{P(E|A)P(A)}{P(E|A)P(A) + P(E|A^c)P(A^c)}$$

This result describes the updating of likelihood of  $A$  from  $P(A)$  to  $P(A|E)$  once  $E$  has been observed. Thomas Bayes actually proved a continuous version of the result, namely, given to  $x$  and  $y$  with conditional distribution  $f(x|y)$  and marginal  $g(y)$ , the conditional distribution of  $y$  given  $x$  is

$$g(y|x) = \frac{f(x|y)g(y)}{\int f(x|y)g(y)dy}$$

While this inversion is quite natural from a probabilistic point of view; Bayes and Laplace went further and considered that the uncertainty of the parameter  $\theta$  of a model could be modelled through a probabilistic distribution  $\pi$  on  $\Theta$  (parameter space) called the prior distribution.

The inference then based on the distribution of  $\theta$  conditional on  $x$ , say,  $\pi(\theta|x)$ , called the posterior distribution, given by

$$\pi(\theta|x) = \frac{f(x|\theta) \pi(\theta)}{\int f(x|\theta) \pi(\theta) d\theta}.$$

The distribution  $\pi(\theta)$  is called a prior distribution on a prior because it quantifies one's uncertainty about  $\theta$ , prior to observing data. The prior may represent the Bayesian's subjective belief/knowledge, in which case it would be a subjective prior.

Alternatively, it could be a conventional prior supposed to represent small or no information.

Recall  $\pi(\theta|x) = \frac{\pi(\theta) f(x|\theta)}{\int \pi(\theta) f(x|\theta) d\theta}$  ----- (\*)

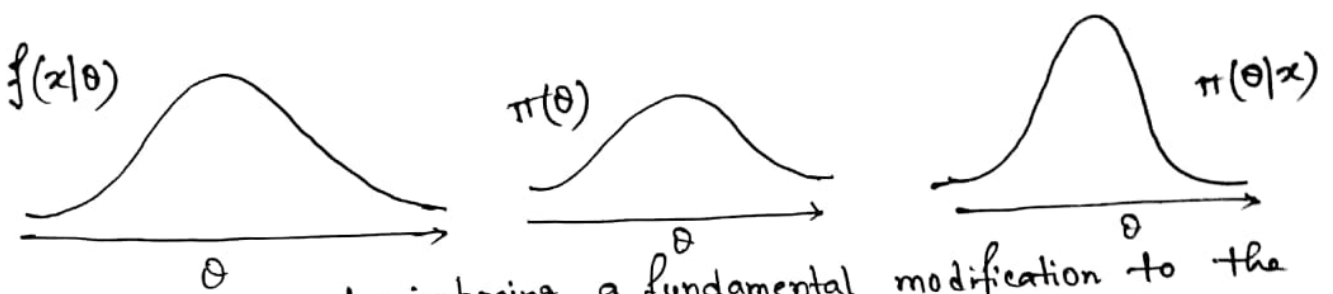
i.e.,  $\pi(\theta|x) \propto \pi(\theta) f(x|\theta)$ .

The numerator in (\*) is the joint density of  $\theta$  and  $x$ .

The denominator  $m(x) \stackrel{\text{def}}{=} \int \pi(\theta) f(x|\theta) d\theta$  is the marginal density of  $x$ .

The posterior distribution combines the prior beliefs about  $\theta$  with the information about  $\theta$  contained in the data  $x$ , to give a composite picture of the final beliefs about  $\theta$ .

Just as prior distribution reflects beliefs about  $\theta$  prior to analysing data,  $\pi(\theta|x)$  reflects the updated beliefs about  $\theta$  after (posterior to) observing data  $x$ .



In summary, by imposing a fundamental modification to the perception of random phenomena, in the sense that parameters diverting random phenomena can also be perceived as "random variables". Bayes and Laplace created "modern" statistical analysis. The use of prior is the best way to summarize the available information (or even lack of information) about  $\theta$  as well as residual uncertainty, thus allowing for incorporation of this imperfect information in the decision process.

Remark:- The posterior distribution contains all the current information about  $\theta$ . Ideally, one might report the entire posterior distribution  $\pi(\theta|x)$  or one could report descriptive summary measures associated with the posterior distribution.

Example:- For a real valued parameter  $\theta$ , one could report the parameter mean  $E(\theta|x) = \int \theta \pi(\theta|x) d\theta$  and the

(H)

posterior variance 
$$\text{Var}(\theta|x) = E \left[ (\theta - E(\theta|x))^2 | x \right]$$

$$= \int (\theta - E(\theta|x))^2 \pi(\theta|x) d\theta.$$

Finally one could use  $\pi(\theta|x)$  to answer question related to more structured problem like estimation or testing.



**Point Estimation:-** Commonly used summaries of location - like mean, median or mode of  $\pi(\theta|x)$  may be used as point estimates of  $\theta$ . Together with a point estimate of  $\theta$ , one also reports a measure of extent of uncertainty associated with the estimate.

**Example:** Together with the posterior mean one also reports posterior variance.

Optimum estimates can be obtained for given loss function. Let,  $L(\theta, a)$  be the loss in estimating  $\theta$  by  $a$ . Bayes estimate of  $\theta$  is obtained by minimizing the average loss

$$E^{\pi(\cdot|x)} L(\theta, a) = \int L(\theta, a) \pi(\theta|x) d\theta \quad \text{w.r.t. } a.$$

The minimizer is called the Bayes estimate of  $\theta$ .

For  $L(\theta, a) = (\theta - a)^2$ , then Bayes estimate is posterior mean,

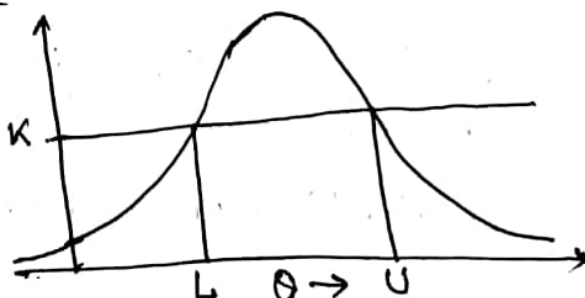
For  $L(\theta, a) = |\theta - a|$ , then Bayes estimate is posterior median.

If one reports posterior median, it is convenient to report a couple of other posterior quantities to give an idea of the posterior variability of  $\theta$ . One may report the first and third posterior quantiles.

**Confidence Set:-** The Bayesian analysis of classical confidence set is called a credible set.

**Definition:-** For  $0 < \alpha < 1$ , a  $100(1-\alpha)\%$  credible set for  $\theta$  is a subset  $C$  of the parameter space  $(H)$  w.r.t.

$$P(\theta \in C|x) = \int_{\theta \in C} \pi(\theta|x) d\theta \geq 1 - \alpha.$$



$$HPD: \{ \theta : \pi(\theta|x) \geq k \}$$

Where as the classical confidence statement don't apply to whether for a given set for a given  $x$  covers the "true"  $\theta$ , this is not the case with credible set does answer a layman's question whether the given set covers the "true"  $\theta$  with probability  $1-\alpha$ . This is because, in the Bayesian approach, "true"  $\theta$  is a random variable, with a data dependent probability distribution (after data has been observed), namely, the posterior distr.

In choosing a credible set for  $\theta$ , it is usually desirable to minimize its size. To this, one should include in the set only those points with the largest posterior density, i.e., the "most likely" values of  $\theta$ .

Definition:- The  $100(1-\alpha)\%$  Highest Posterior Density <sup>(HPD)</sup> credible set for  $\theta$  is the subset  $C$  of  $\Theta$  of the form

$$C = \{ \theta \in \Theta : \pi(\theta|x) \geq k(\alpha) \}, \text{ where}$$

$k(\alpha)$  is the largest constant respect to  $P(\theta \in C | x) \geq 1-\alpha$ .

Example of credible Sets:-  $x_1, \dots, x_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ ;  $\sigma^2$  known.

Let  $\theta \sim N(\mu, \tau^2)$  as before ( $\mu, \tau^2$ ) is known, then highest posterior density (HPD) credible set for  $\theta$  is an equal tailed interval  $C$  centred at the posterior mean

$$C = E(\theta|x) \pm z_{\alpha/2} \sqrt{\text{Var}(\theta|x)}$$

Example:-  $x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \text{Cauchy}(\theta)$ ;  $\theta > 0$

$$f(x|\theta) = \frac{1}{\pi(1+(x-\theta)^2)} ; -\infty < x < \infty, \text{ where } \theta \text{ is the location parameter.}$$

A reasonable objective prior one can easily find a  $100(1-\alpha)\%$  HPD credible set using computer numerically based on posterior densities

$$\pi(\theta|x) \propto \prod_{i=1}^n [1 + (\theta - x_i)^2]^{-1} ; \theta > 0.$$

In contrast, it is not at all clear how to construct a classical confidence set for this problem for fixed sample size situation. For very large sample sizes, we may use the asymptotic normality of sample median to produce an asymptotically correct frequentist CI for  $\theta$ .

Hypothesis Testing:- It is possible to derive posterior probability of a hypothesis  $H_0$ , i.e.,  $P^\pi(H_0 | \underline{x})$ .  
 Suppose we want to test  $H_0: \theta \leq \theta_0$  Vs.  $H_1: \theta > \theta_0$ .

One calculates

$$P(H_0 | \underline{x}) = P(\theta \leq \theta_0 | \underline{x}) = \int_{\theta \leq \theta_0} \pi(\theta | \underline{x}) d\theta$$

Exercise:- Consider the problem of estimation of a real parameter  $\theta$  with the loss function:

$$L(\theta, a) = \begin{cases} k_0(\theta - a) & \text{if } \theta - a \geq 0 \\ k_1(a - \theta) & \text{if } \theta - a < 0 \end{cases}$$

Show that the Bayes estimate is given by the quantity of order  $\frac{k_0}{k_0 + k_1}$  of the posterior distribution (assume, for simplicity uniqueness of the quantile).

Example 1: Normal Means:-

Hypothesis Testing:- Typically for tests like  $H_0: \theta \in \mathcal{H}_0$  Vs.  $H_1: \theta \in \mathcal{H} - \mathcal{H}_0$ , one derives the testing rule using the posterior probability  $P^\pi(H_0 | \underline{x})$  on  $H_0$  being true.

Example: 1. Normal Means:-

Suppose given  $\theta$ ;  $X_1, X_2, \dots, X_n$  are iid  $N(\theta, \sigma^2)$  where  $\sigma^2$  is the known parameter and the parameter of interest is  $\theta$ ; say,  $\theta \in \mathbb{R}$ .

A convenient and desirable prior distr. for  $\theta$  is a normal distr. with suitable prior mean and variance, say,  $\theta \sim \pi(\theta) = N(\mu, \tau^2)$  where,  $\mu$  and  $\tau^2$  are both known. The prior variance  $\tau^2$  is a measure of strength of our belief in the prior mean in the sense that larger the  $\tau^2$  the less sure we are about our prior guess. The likelihood is —

$$L(\theta) = f(x_1, \dots, x_n | \theta) = \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2\right]$$

Prior density,  $\pi(\theta) = \frac{1}{\tau\sqrt{2\pi}} \exp\left[-\frac{1}{2\tau^2}(\theta-\mu)^2\right]$

Posterior density,  $\pi(\theta|\underline{x}) = \frac{f(\underline{x}|\theta)\pi(\theta)}{\int f(\underline{x}|\theta)\pi(\theta)d\theta}$

$\propto f(\underline{x}|\theta)\pi(\theta) = \exp\left(-\frac{M}{2}\right)$ , where,

$$M = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \theta)^2 + \frac{1}{\tau^2} (\theta - \mu)^2$$

$$= \frac{1}{\sigma^2} \left( \sum x_i^2 - 2\theta \sum x_i + n\theta^2 \right) + \frac{1}{\tau^2} (\theta^2 - 2\theta\mu + \mu^2)$$

$$= \theta^2 \left( \frac{1}{\tau^2} + \frac{n}{\sigma^2} \right) - 2\theta \left( \frac{\mu}{\tau} + \frac{\sum x_i}{\sigma^2} \right) + \frac{\mu^2}{\tau^2} + \frac{\sum x_i^2}{\sigma^2}$$

$$= \theta^2 a - 2b\theta + c$$

$$= a \left( \theta - \frac{b}{a} \right)^2 + c - \frac{b^2}{a}$$

Thus,  $\pi(\theta|\underline{x}) \propto \exp\left[-\frac{1}{2} a \left( \theta - \frac{b}{a} \right)^2\right]$

$$\therefore \theta|\underline{x} \sim N\left(\frac{b}{a}, \frac{1}{a}\right) = N\left(\frac{\frac{\sum x_i^2}{\sigma^2} + \frac{\mu}{\tau^2}}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}, \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}\right)$$

$E(\theta|\underline{x}) = \frac{\frac{\mu}{\tau^2} + \frac{n\bar{x}}{\sigma^2}}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}$  and  $\text{Var}^\pi(\theta|\underline{x}) = \frac{1}{\frac{1}{\tau^2} + \frac{n}{\sigma^2}}$ .

(posterior mean) (posterior variance)

If precision is measured by reciprocal of variance, posterior precision is  $\frac{1}{\tau^2} + \frac{n}{\sigma^2}$ . In the layout of the data, the precision in the prior knowledge has increased and is reflected in the increased value of the posterior precision.

Note that, the posterior mean  $\theta$  is an weighted average of the prior estimate  $\mu$  and the classical estimate  $\bar{x}$  with weights being proportional to the corresponding precisions

$$E(\theta|\underline{x}) = \left( \frac{1/\tau^2}{1/\tau^2 + n/\sigma^2} \right) \mu + \left( \frac{n/\sigma^2}{1/\tau^2 + n/\sigma^2} \right) \bar{x}$$

It may be considered as a modification of the classical estimate  $\bar{x}$  in the light of prior information about  $\theta$  or as a modification of the prior estimate  $\mu$  of  $\theta$  in the light of the data and note that as  $n \rightarrow \infty$ ;  $E(\theta | \bar{x}) \rightarrow \bar{x}$  as  $n \rightarrow \infty$ .

Thus if the prior information is small or there are lots of data points the posterior mean is close to the classical estimate  $\bar{x}$ .

In both cases  $\pi(\theta | \bar{x}) \approx N(\bar{x}, \sigma^2/n)$ ; as  $n \rightarrow \infty$ ,  $v(\theta | \bar{x}) \rightarrow 0$ .

Thus as  $n \rightarrow \infty$  the posterior distr. becomes more and more concentrated around  $\bar{x}$ .

- The posterior depends both on prior and  $\bar{x}$ .
- As amount of data increases the influence of data leads to wash away the prior.

Example: 2. Suppose  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$ ; let  $p$  have a prior distr.  $\pi(p)$ .

We will consider a family of priors for  $p$  which simplifies the calculation of posterior and then consider some commonly used priors from this family.

Let  $\pi(p) \propto p^{\alpha-1} (1-p)^{\beta-1}$ ; Beta( $\alpha, \beta$ );  $\alpha > 0, \beta > 0$ .

Prior mean =  $\frac{\alpha}{\alpha+\beta}$ ; Variance =  $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$ .

The posterior,  $\pi(p | \bar{x}) \propto p^{\alpha+\gamma-1} (1-p)^{\beta+(n-\gamma)-1}$ ; where  $\gamma = \sum_{i=1}^n X_i$ .

The posterior distr. of  $p$  is Beta( $\alpha + \sum X_i, \beta + n - \sum X_i$ ).

$$E(p | \bar{x}) = \frac{\alpha + \sum X_i}{\alpha + \beta + n}$$

$$\text{Var}(p | \bar{x}) = \frac{(\alpha + \sum X_i)(\beta + n - \sum X_i)}{(\alpha + \beta + n)^2 (\alpha + \beta + n + 1)}$$

If  $n$  is large then  $E(p | \bar{x}) \approx \hat{p} = \frac{\sum X_i}{n}$  is the classical estimate and  $\text{Var}(p | \bar{x})$  is also quite small. The posterior distr. is concentrated around  $\hat{p} = \sum X_i / n$ .

We can interpret this as an illustration of a fact mentioned above: when we have lots of data, the data tend to wash away the influence of prior. The posterior mean can be written as a weighted average of the prior mean (prior estimate) and the classical estimate (MLE)  $\hat{p} = \frac{\sum X_i}{n}$ .

Note that the posterior mean  $E(\theta|X)$  can be written as

$$\frac{\alpha + \beta}{\alpha + \beta + n} \left( \frac{\alpha}{\alpha + \beta} \right) + \frac{n}{\alpha + \beta + n} \left( \frac{\sum X_i}{n} \right).$$

- As in example 1, we see that the prior mean and the classical estimate are being combined in a convex way.

Taking  $\alpha = 1 = \beta$ , for the uniform distr., the posterior mean is  $\frac{\sum X_i + 1}{n + 2}$ .  
 The prior was used by Bayes and Laplace.  
 (Uniform prior)

- If  $\alpha = \beta = 1/2$ , we have what is called the Jeffery's prior.  
 Jeffery's prior  $\pi(\theta) \propto \sqrt{\det(I(\theta))}$ ; where  $I(\theta)$  is the fisher information based on single observation.

- The posterior mean corresponding to Jeffery's prior is given by  $\frac{\sum X_i + 1/2}{n + 1}$ . Jeffery's prior is a very widely used and popular prior in the case of one dimensional parameter.

- When  $\alpha = \beta = 0$ , the posterior mean is equal to MLE =  $\frac{\sum X_i}{n}$ .

Remark: If we had wished to test  $H_0: \mu = \mu_0$  vs.  $H_1: \mu \neq \mu_0$  or  $H_1: \mu > \mu_0$ .  
 We have to choose the prior in a careful way since any continuous prior on  $\mu$  would give zero probability to  $H_0$ . In these problems the bayesian answers tend to differ substantially from classical answers.

$$\begin{aligned} X_1, \dots, X_n & \stackrel{iid}{\sim} c(\theta, 1) \\ X_i - \theta & \stackrel{iid}{\sim} c(0, 1) \\ \bar{X} - \theta & \stackrel{iid}{\sim} c(0, 1) \end{aligned}$$

Improper Priors:- Consider the "normal mean" example:

$X_1, \dots, X_n \stackrel{iid}{\sim} N(\theta, \sigma^2)$ ,  $\sigma^2$  is known,  $\theta \sim N(\mu, \tau^2)$ .

As  $\tau^2$  becomes larger, the prior becomes more and more diffuse/flat. If  $\tau^2$  is very large ( $\tau^2 \rightarrow \infty$ ),

$$\pi(\theta | \underline{x}) \approx N(\bar{x}, \sigma^2/n) \text{ density for } \theta,$$

Putting this in another way, if we let  $\tau^2 \rightarrow \infty$ , then in the limiting case it would give an improper uniform prior over the whole real line

$$\pi(\theta) = c, \quad -\infty < \theta < \infty \quad (c \text{ is an arbitrary constant})$$

which, by the Bayes formula, will give a  $N(\bar{x}, \frac{\sigma^2}{n})$  posterior for  $\theta$ .

Note that  $\int \pi(\theta) d\theta = \infty$ . Such a prior distn is called improper. We call a prior density  $\pi(\theta)$  proper if it integrates to 1.

If  $\int \pi(\theta) d\theta$  is finite, it is called an unnormalised density and can be normalized — multiplied by a constant — to integrate to 1.

In the above example, despite the impropriety of the prior distn, the posterior distn is proper given just one observation.

With an improper prior, we can proceed with the Bayes formula and define an unnormalized posterior density by

$$\pi(\theta | \underline{x}) \propto f(\underline{x} | \theta) \pi(\theta).$$

If  $m(\underline{x}) = \int f(\underline{x} | \theta) \pi(\theta) d\theta < \infty$  this leads to a proper posterior density. Bayesian methods apply as long as the posterior distribution is proper.

## Noninformative Priors (Objective Priors):-

If no prior information is available, or, it is very difficult to quantify the available information into a prior distribution, one can still use Bayesian methods using what are called noninformative (objective) priors.

Standard noninformative priors are:

Uniform prior :  $\pi(\theta) = \text{constant}$

Jeffreys prior :  $\pi(\theta) = \sqrt{\det(I(\theta))}$ , where  $I(\theta)$  is Fisher Information matrix.

Reference priors (Bernardo, Berger)

Several methods for construction of noninformative priors are available.

Standard noninformative priors for location parameter ( $\mu$ ) and scale parameter ( $\sigma$ ) are:  $\pi(\mu) = c$ ,  $-\infty < \mu < \infty$  and  $\pi(\sigma) \propto \frac{1}{\sigma}$ ,  $0 < \sigma < \infty$ .

$$f(x, \theta) = f(x - \theta), \quad f(x, \sigma) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}\right), \quad \sigma > 0.$$

Note: Objective priors are usually improper. To be usable they must have proper posteriors. The uniform, Jeffreys' prior and reference priors are examples of objective priors. All of them produce a posterior mean that is close to the MLE mean for small  $n$ .



Conjugate Priors:- In the Normal and Bernoulli examples, we saw that by taking normal and Beta as priors respectively the posterior remains within the normal and beta families. The property that posterior distr. is of the same parametric form as the prior distr. is called conjugacy. The corresponding priors are called conjugate priors. These priors also have great advantage in terms of computational ease and can also be interpreted as "additional data".

Example: ①  $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\theta); \theta > 0.$

$$\pi(\theta) = \frac{\beta^\alpha \theta^{\alpha-1} e^{-\beta\theta}}{\Gamma(\alpha) \theta^{\sum X_i + \alpha - 1} e^{-(\beta+n)\theta}} \quad [\text{Gamma}(\alpha, \beta)]$$

Posterior density i.e.,  $\text{Gamma}(\alpha + \sum X_i, \beta + n)$

$$E(\theta | \underline{x}) = \frac{\alpha + \sum X_i}{\beta + n}; \quad \text{Var}(\theta | \underline{x}) = \frac{\alpha + \sum X_i}{(\beta + n)^2}$$

②  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2); \mu \text{ known, } \sigma^2 \text{ unknown.}$   
Likelihood is  $f(\underline{x} | \sigma^2) \propto (\sigma^2)^{-n/2} \exp\left[-\frac{n\upsilon}{2\sigma^2}\right];$

$$\text{where } \upsilon = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

The conjugate prior density is the inverse gamma.

$$\pi(\sigma^2) \propto (\sigma^2)^{-(\alpha+1)} e^{-\beta/\sigma^2}$$

$$\text{Posterior: } \pi(\sigma^2 | \underline{x}) \propto (\sigma^2)^{-(\alpha+1+n/2)} \exp\left[-\frac{1}{\sigma^2} \left(\beta + \frac{n\upsilon}{2}\right)\right]$$

$$L(\sigma^2, a) = \frac{(a - \sigma^2)^2}{\sigma^4} = \left(\frac{a}{\sigma^2} - 1\right)^2$$

To minimize w.r.t.  $a$ ,

$$\int L(\sigma^2, a) \pi(\sigma^2 | \underline{x}) d\sigma^2$$

$$= \int (a - \sigma^2)^2 (\sigma^2)^{-2} \times \left\{ \text{density of IG}\left(\alpha + \frac{n}{2}, \beta + \frac{n\upsilon}{2}\right) \text{ for } \sigma^2 \right\} d\sigma^2$$

$$= \text{constant} \times \int (a - \sigma^2)^2 \times \left\{ \text{density of IG}\left(\alpha + \frac{n}{2} + 2, \beta + \frac{n\upsilon}{2}\right) \text{ for } \sigma^2 \right\} d\sigma^2$$

$E \{ L(\theta^2, a) | X \}$  is minimized at

$$a = \text{Mean of IG} \left( \alpha + 2 + \frac{n}{2}, \beta + \frac{nV}{2} \right)$$

$$= \frac{\beta + \frac{nV}{2}}{\alpha + 1 + \frac{n}{2}} = \text{mode of posterior distr.}$$

[ If  $\theta \sim \text{IG}(\alpha, \beta)$  ;  $E(\theta) = \frac{\beta}{\alpha - 1}$  ;  $\alpha > 1$

$$V(\theta) = \frac{\beta^2}{(\alpha - 1)^2 (\alpha - 2)} ; \alpha > 2$$

$$\text{Mode}(\theta) = \frac{\beta}{\alpha + 1} . ]$$

③ Exponential Family:-  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x|\theta)$

$$f(x|\theta) = \exp \{ \theta t(x) + c(\theta) \} h(x) ; \theta \in (a, b) ; a, b \text{ could be } \pm \infty.$$

Assume all regularity conditions and smoothness of  $c(\theta)$  ; consider the problem of estimating  $\mu = E_{\theta}(t(x))$ .

Note:-  $\mu = \mu(\theta) = -c'(\theta)$

$$E_{\theta} \left[ \frac{\partial}{\partial \theta} \log f(x|\theta) \right] = 0 = E_{\theta} \left[ c'(\theta) + t(x) \right]$$

$$\Rightarrow E_{\theta}(t(x)) = -c'(\theta).$$

Now,  $0 < I(\theta) = \text{Var} \left[ \frac{\partial}{\partial \theta} \log f(x|\theta) \right]$

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

$$= -c''(\theta)$$

$$\text{i.e., } \frac{d\mu}{d\theta} = -c''(\theta) > 0$$

implying  $\mu$  is a strictly increasing function of  $\theta$ .

MLE of  $\mu$  :

$$L = f(x_1, \dots, x_n | \theta) \times e^{nc(\theta) + \theta T} ; \text{ where, } T = \sum_{i=1}^n t(x_i).$$

$$\frac{\partial \log L}{\partial \theta} = 0 \Rightarrow nc'(\theta) + T = 0$$

$$\Rightarrow -c'(\theta) = T/n$$

$$\Rightarrow \hat{\mu}_{\text{MLE}} = T/n.$$

Conjugate Prior:

$\pi(\theta) \propto \exp[mc(\theta) + \theta s]$ ;  $\theta \in (a, b)$ . We assume  $\pi(a) = \pi(b) = 0$ .

Posterior distribution:-

$$\pi(\theta | \underline{x}) \propto e^{(m+n)c(\theta) + \theta(s+T)}$$

[Another way of looking at conjugate prior is as follows:  
Consider a prior  $\pi(\theta) = 1$ . Take  $m$  to be +ve integer and think of a hypothetical sample of size  $m$  with hypothetical data  $x'_1, \dots, x'_m$  s.t.  $\sum_{i=1}^m t(x'_i) = s$  then

$$\begin{aligned} \pi(\theta | x'_1, \dots, x'_m) &\propto \int (x'_1, \dots, x'_m | \theta) \pi(\theta) \\ &= \exp\left\{\theta \cdot \sum_{i=1}^m t(x'_i) + mc(\theta)\right\} \\ &= e^{\theta s + mc(\theta)}. \end{aligned}$$

Prior mean:

$$\begin{aligned} E^{\pi}(\mu) &= E^{\pi}(-c'(\theta)) \\ &= -k \int_a^b c'(\theta) e^{mc(\theta) + \theta s} d\theta \\ &= -\frac{k}{3k} \int_a^b \frac{d}{d\theta} (e^{mc(\theta)}) \cdot e^{\theta s} d\theta \\ &= -\frac{k}{3k} \left[ \left\{ e^{\theta s} e^{mc(\theta)} \right\}_a^b - \int_a^b s e^{\theta s} e^{mc(\theta)} d\theta \right] \\ &= -\frac{1}{m} \left[ \int_a^b \pi(\theta) d\theta - s \int_a^b \pi(\theta) d\theta \right] \\ &= \frac{s}{m}. \end{aligned}$$

Similar calculations show that posterior mean of  $\mu$  is

$$\frac{s+T}{m+n} = \frac{m}{m+n} \underbrace{\left(\frac{s}{m}\right)}_{\text{prior guess}} + \frac{n}{m+n} \underbrace{\left(\frac{T}{n}\right)}_{\text{classical estimate}}$$

So, the posterior mean is a weighted average of the prior guess  $\frac{s}{m}$  and MLE  $\frac{T}{n}$  and the weights are proportional to the precision parameter  $m$  and sample size  $n$ .

For multiparameter exponential families given by

$$p(x|\underline{\theta}) = \exp \left\{ c(\underline{\theta}) + \sum_{i=1}^d \theta_i t_i(x) + h(x) \right\}$$

The conjugate prior now takes the form

$$p(\underline{\theta}) = \exp \left\{ mA(\underline{\theta}) + \sum_{i=1}^d \theta_i s_i \right\}$$

④.  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\gamma, \theta)$ ;  $\gamma$  is known ( $> 0$ ).

Likelihood function:

$$f(x|\theta) = \prod_{i=1}^n \frac{\theta^\gamma}{\Gamma(\gamma)} x^{\gamma-1} e^{-\theta x_i} = \frac{\theta^{n\gamma}}{(\Gamma(\gamma))^n} \left( \prod_{i=1}^n x_i \right)^{\gamma-1} e^{-\theta \sum_{i=1}^n x_i}$$

conjugate prior:  $\pi(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}$ ;  $\theta > 0$

Posterior density:  $\pi(\theta|x) \propto \theta^{\alpha+n\gamma-1} e^{-(\beta+\sum x_i)\theta}$ ,  
a  $\text{Gamma}(\alpha+n\gamma, \beta+\sum_{i=1}^n x_i)$  density.

Posterior mean and variance:

$$E(\theta|x) = \frac{\alpha+n\gamma}{\beta+\sum x_i}, \quad \text{Var}(\theta|x) = \frac{\alpha+n\gamma}{(\beta+\sum x_i)^2}$$

⑤  $f(x|\theta)$  is  $N(\theta, 1)$  density.

$$f(x|\theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(x-\theta)^2\right) \\ = \exp\left(\theta x - \frac{1}{2}\theta^2\right) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right)$$

So,  $c(\theta) = -\frac{1}{2}\theta^2$

$t(x) = x$

$T = \sum_{i=1}^n t(x_i) = n\bar{x}$

$\mu = E_{\theta}(x) = \theta = -c'(\theta)$ .

$\pi(\theta) \propto \exp\left[-\frac{m}{2}\theta^2 + \theta s\right]$

$\propto \exp\left[-\frac{m}{2}\left(\theta - \frac{s}{m}\right)^2\right]$  which is  $N\left(\frac{s}{m}, \frac{1}{m}\right)$ .

Prior mean:  $\frac{s}{m}$ .

Posterior mean:  $\frac{s+n\bar{x}}{m+n}$ .

Predictive distribution:- Distribution of a future observation.  
 Let  $x_1, \dots, x_n$  be i.i.d.  $f(x|\theta)$  and  $\theta \sim \pi(\theta)$ .

Prior predictive distribution:- Prior predictive distribution of a future observation  $X$  is given by density  $m(x) = \int f(x|\theta)\pi(\theta)d\theta$ .

Intuitively  $\theta$  is known,  $X$  will follow  $f(x|\theta)$ ,  $\textcircled{H}$  we don't know  $\theta$  but  $\theta \sim \pi(\theta)$ . Therefore, predictive distr. is obtained as an average  $f(x|\theta)$  with respect to  $\pi(\theta)$ .  
 Similarly, we can define prior predictive distr. of first  $n$  observations as

$$m(x_1, x_2, \dots, x_n) = \int f(x_1, \dots, x_n | \theta) \pi(\theta) d\theta.$$

Posterior Predictive distribution:- Distribution of  $X_{n+1}$  given  $x_1, \dots, x_n$   
 is  $f(x_{n+1} | x_1, \dots, x_n) = \int f(x_{n+1} | \theta) \pi(\theta | x_1, \dots, x_n) d\theta$ .

where  $\pi(\theta | x_1, \dots, x_n)$  is density of the posterior distr. of  $\theta$  given  $x_1, \dots, x_n$ . Then  $f(x_{n+1} | x_1, \dots, x_n)$  is called posterior predictive distr. of  $X_{n+1}$  given  $x_1, \dots, x_n = (x_1, \dots, x_n)$ .

Example:- 1.  $X_1, \dots, X_n \sim \text{Ber}(\theta)$

$$f(x|\theta) = \theta^{\sum x_i} (1-\theta)^{n - \sum x_i}$$

Consider  $\pi(\theta) = 1$ , then

$$f(x_{n+1} | x_1, \dots, x_n) = \int_0^1 \theta^{\sum x_i + x_{n+1}} (1-\theta)^{n - \sum x_i - x_{n+1}} \pi(\theta | x) d\theta$$

$$= \frac{\int_0^1 \theta^{\sum x_i + x_{n+1}} (1-\theta)^{n - \sum x_i - x_{n+1}} d\theta}{\int_0^1 \theta^{\sum x_i} (1-\theta)^{n - \sum x_i} d\theta}$$

$$= \frac{\text{Beta}(\sum x_i + x_{n+1} + 1, n - \sum x_i - x_{n+1} + 1)}{\text{Beta}(\sum x_i + 1, n - \sum x_i + 1)} = \frac{n - \sum x_i + 1}{n + 2}$$

$$P(X_{n+1} = 1 | X_1 = 1, \dots, X_n = 1) = \frac{n+1}{n+2}$$

### Example 2: (Normal Means)

$X_1, X_2, \dots, X_n$  be iid  $N(\theta, \sigma^2)$ ,  $\sigma^2$  known.

$\theta \sim \pi(\cdot) = N(\mu, \tau^2)$ ,  $\mu, \tau^2$  known.

$\pi(\theta | \underline{x}) = \text{Normal with mean} = \frac{n\bar{x}}{\sigma^2} + \frac{\mu}{\tau^2}$  and

$$\text{variance} = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}}.$$

$$f(x_{n+1} | x_1, \dots, x_n) = \int N(x_{n+1} | \theta, \sigma^2) \pi(\theta | \underline{x}) d\theta$$

which is a normal distr. with

$$\text{mean} = E[X_{n+1} | \underline{x}] = E[E(X_{n+1} | \underline{x}, \theta)]$$
$$= E[\theta | \underline{x}]$$

$$\text{Var}(X_{n+1} | \underline{x}) = \text{Var}[E\{X_{n+1} | \theta, \underline{x}\}]$$
$$+ E[\text{Var}\{X_{n+1} | \theta, \underline{x}\}]$$

$$= \text{Var}(\theta | \underline{x}) + \sigma^2$$

$$= \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\tau^2}} + \sigma^2.$$

Exercise:- For the above example, show that the prior predictive distr. (joint distr.) of  $X_1, X_2, \dots, X_n$  is normal with  $E(X_i) = \mu$  and  $\text{Var}(X_i) = \tau^2 + \sigma^2 \forall i$  and  $\text{Cov}(X_i, X_j) = \tau^2 \forall i \neq j$ .  
Prediction of a future observation by a single number  $t(x_1, x_2, \dots, x_n)$  based on  $x_1, x_2, \dots, x_n$  with squared error loss amounts to considering prediction loss  $E[(X_{n+1} - t)^2 | x_1, \dots, x_n]$  which is minimum at  $t = E(X_{n+1} | x_1, \dots, x_n)$ .

Multiparameter models :-

Ex. 1. (Normal)

$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ; both  $\mu, \sigma^2$  unknown.

A noninformative prior for  $\mu$  and  $\sigma^2$  is

$\pi(\mu, \sigma^2) \propto \frac{1}{\sigma^2}$  (conventional improper prior uniform on  $(\mu, \log \sigma^2)$ )

Joint posterior density

$\pi(\mu, \sigma^2 | \underline{x}) \propto (\sigma^2)^{-n/2-1} \exp\left[-\frac{1}{2\sigma^2} \{(n-1)s^2 + n(\mu - \bar{x})^2\}\right]$ ,

where  $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ .

Conditional posterior density of  $\mu$  given  $(\sigma^2, \underline{x})$  is given by

$\pi(\mu | \sigma^2, \underline{x}) \sim N(\bar{x}, \frac{\sigma^2}{n})$ .

Marginal posterior distn:

$\pi(\sigma^2 | \underline{x}) \propto \int (\sigma^2)^{-n/2-1} \exp\left[-\frac{1}{2\sigma^2} \{(n-1)s^2 + n(\mu - \bar{x})^2\}\right] d\mu$   
 $\propto (\sigma^2)^{-\frac{(n+1)}{2}} \cdot \exp\left[-\frac{1}{2\sigma^2} (n-1)s^2\right]$

This is a scale inverse  $\chi^2$  density.

$\pi(t | \underline{x}) \propto \left[\frac{(n-1)s^2}{t}\right]^{-\frac{(n+1)}{2}} \exp\left[-\frac{t}{2}\right] \frac{(n-1)s^2}{t^2}$

$\propto t^{\frac{n-1}{2}-1} \exp\left[-\frac{t}{2}\right]$

Conditional distn of  $\frac{(n-1)s^2}{\sigma^2}$  given  $\underline{x}$  is  $\chi^2_{n-1}$ .

[ Sampling theory result :-

Distn of  $\frac{(n-1)s^2}{\sigma^2}$  (given  $\mu, \sigma^2$ ) is  $\chi^2_{n-1}$ .

Marginal Posterior distribution of  $\mu$ :

$$\pi(\mu | \underline{x}) = \int_0^{\infty} \pi(\mu, \sigma^2 | \underline{x}) d\sigma^2$$

$$\propto \int_0^{\infty} (\sigma^2)^{-\eta/2-1} e^{-\frac{A}{2\sigma^2}} d\sigma^2$$

$$\left[ \begin{aligned} A &= (n-1)s^2 + n(\mu - \bar{x})^2 \\ z &= \frac{A}{2\sigma^2}, \quad \sigma^2 = \frac{A}{2z} \end{aligned} \right]$$

$$= \int_0^{\infty} \left(\frac{A}{2z}\right)^{-\eta/2-1} e^{-z} \frac{A}{2z^2} dz$$

$$\propto A^{-\eta/2} \int_0^{\infty} z^{\eta/2-1} e^{-z} dz$$

$$\propto \left[ 1 + \frac{n(\mu - \bar{x})^2}{(n-1)s^2} \right]^{-\eta/2}$$

This is the  $t$ -distribution  $t_{n-1}(\bar{x}, \frac{s^2}{n})$ .

[  $\theta \sim t_{\nu}(\mu, \sigma^2)$  mean  $p(\theta) = \text{constant}$  ]

Standard  $t \sim t_{\nu}(0, 1)$  denoted by  $t_{\nu}$ .

Note that  $\frac{\mu - \bar{x}}{s/\sqrt{n}}$  (given  $\underline{x}$ )  $\sim t_{n-1}$ .

[ Sampling theory result:  $\frac{\mu - \bar{x}}{s/\sqrt{n}}$  (given  $\mu, \sigma^2$ )  $\sim t_{n-1}$  ]



## Multiparameter Models :-

Example (Normal Data with a conjugate prior distribution) :-

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ;  $\mu, \sigma^2$  both unknown.

Prior:  $\pi(\mu, \sigma^2) = \pi(\mu | \sigma^2) \pi(\sigma^2)$

$$\mu | \sigma^2 \sim \pi(\mu | \sigma^2) = N\left(\mu_0, \frac{\sigma^2}{k_0}\right)$$

$$\sigma^2 \sim \pi(\sigma^2) = \text{Inverse-}\chi^2(\nu_0, \sigma_0^2) \quad \left[ \text{d.f.} = \nu_0, \text{scale} = \sigma_0^2 \right]$$

Then  $\pi(\mu, \sigma^2) \propto \frac{1}{(\sigma^2)^{1/2}} \exp\left[-\frac{k_0}{2\sigma^2} (\mu - \mu_0)^2\right]$

$$\times (\sigma^2)^{-\left(\frac{\nu_0}{2} + 1\right)} e^{-\nu_0 \sigma_0^2 / 2\sigma^2}$$

We call this the  $N\text{-Inv-}\chi^2\left(\mu_0, \frac{\sigma^2}{k_0}; \nu_0, \sigma_0^2\right)$  density.

## Joint Posterior Density :-

$$\pi(\mu, \sigma^2 | \underline{x}) \propto (\sigma^2)^{-\frac{(\nu_0+3)}{2}} \exp\left[-\frac{1}{2\sigma^2} \left\{ \nu_0 \sigma_0^2 + k_0 (\mu - \mu_0)^2 \right\}\right]$$

$$\times \frac{1}{(\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2} \left\{ (n-1)s^2 + n(\mu - \bar{x})^2 \right\}\right]$$

$$= N\text{-Inv-}\chi^2\left(\mu_n + \frac{\sigma^2}{k_n}; \nu_n, \sigma_n^2\right)$$

where,  $\mu_n = \frac{k_0}{k_0+n} \mu_0 + \frac{n\bar{x}}{k_0+n}$ ;  $k_n = k_0 + n$ ;  $\nu_n = \nu_0 + n$ .

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + (n-1)s^2 + \frac{k_0 n}{k_0+n} (\bar{x} - \mu_0)^2 \quad (\text{check})$$

## Conditional Posterior Density :-

$$\pi(\mu_n | \sigma_n^2, \underline{x}) : \mu | \sigma_n^2, \underline{x} \sim N\left(\mu_n, \frac{\sigma_n^2}{k_n}\right)$$

Marginal Posterior Densities :-  $\pi(\sigma_n^2 | \underline{x})$  and  $\pi(\mu | \underline{x})$ :

$\pi(\sigma_n^2 | \underline{x})$  is an  $\text{Inv-}\chi^2(\nu_n, \sigma_n^2)$  density.

$$\pi(\mu | \underline{x}) \propto \left[ 1 + \frac{k_n (\mu - \mu_n)^2}{\nu_n \sigma_n^2} \right]^{-(\nu_n+1)/2} \quad (\text{check})$$

is a  $t$ -distribution with location =  $\mu_n$ , scale =  $\frac{\sigma_n}{\sqrt{k_n}}$  and degree of freedom =  $\nu_n$ .

The Multivariate Normal Model:-

$x_1, x_2, \dots, x_n$  iid  $N_p(\mu, \Sigma)$ ;  $\Sigma$  is p.d.

$$f(x_1, x_2, \dots, x_n | \mu, \Sigma) \propto |\Sigma|^{-n/2} \exp \left[ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)' \Sigma^{-1} (x_i - \mu) \right]$$

Consider  $\Sigma$  is known

Conjugate prior distr. for  $\mu$  is:  $\mu \sim N(\mu_0, \Lambda_0)$ ;  $\mu_0, \Lambda_0$  known.

Posterior density of  $\mu$  is:

$$\pi(\mu | \underline{x}, \Sigma) \propto \exp \left[ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)' \Sigma^{-1} (x_i - \mu) - \frac{1}{2} (\mu - \mu_0)' \Lambda_0^{-1} (\mu - \mu_0) \right]$$

$$\begin{aligned} \text{Now, } & \sum_{i=1}^n (x_i - \mu)' \Sigma^{-1} (x_i - \mu) + (\mu - \mu_0)' \Lambda_0^{-1} (\mu - \mu_0) \\ &= \sum_{i=1}^n (x_i - \bar{x})' \Sigma^{-1} (x_i - \bar{x}) + n(\mu - \bar{x})' \Sigma^{-1} (\mu - \bar{x}) \\ & \quad + (\mu - \mu_0)' \Lambda_0^{-1} (\mu - \mu_0) \\ &= \mu' (n\Sigma^{-1} + \Lambda_0^{-1}) \mu - 2\mu' (n\Sigma^{-1} \bar{x} + \Lambda_0^{-1} \mu_0) + \text{constant (free of } \mu) \\ &= (\mu - \mu_n)' \Lambda_n^{-1} (\mu - \mu_n) \end{aligned}$$

$$\text{where, } \mu_n = (\Lambda_0^{-1} + n\Sigma^{-1})^{-1} (\Lambda_0^{-1} \mu_0 + n\Sigma^{-1} \bar{x})$$

$$\Lambda_n^{-1} = \Lambda_0^{-1} + n\Sigma^{-1}$$

$$\text{Then } \pi(\mu | \underline{x}, \Sigma) \propto \exp \left[ -\frac{1}{2} (\mu - \mu_n)' \Lambda_n^{-1} (\mu - \mu_n) \right]$$

$$= N(\mu_n, \Lambda_n) \text{ density.}$$

The parameter mean  $\mu_n$  is a weighted average of  $\bar{x}$  and the prior mean  $\mu_0$  with weights given by data ( $\bar{x}$ ) and prior precision matrices.

Consider now the case with unknown  $\Sigma$

Conjugate prior distr. for  $(\mu, \Sigma)$  is normal-inverse Wishart  
 $(\mu_0, \frac{\Sigma_0}{k_0}; \nu_0, \Lambda_0)$ .

So,  $\mu | \Sigma \sim N(\mu_0, \frac{\Sigma}{k_0})$ ,  $\Sigma \sim \text{Inverse-Wishart}_{\nu_0}(\Lambda_0^{-1})$ .

If  $W^{-1} \sim \text{Wishart}_{\nu}(S)$  ( $\nu = \text{d.f.}$ ,  $S = p \times p$  scale matrix)  
 Then  $W \sim \text{Inverse-Wishart}_{\nu}(S^{-1})$ .

$$f(W) = \left( 2^{\nu p/2} \prod_{i=1}^p \frac{\Gamma(\frac{\nu+1-i}{2})}{\Gamma(\frac{\nu+1-i}{2})} \right)^{-1}$$

$$\times |S|^{\nu/2} |W|^{-(\nu+p+1)/2} \exp\left[-\frac{1}{2} \text{tr}(SW^{-1})\right]; W \text{ is p.d.}$$

The joint prior density is:

$$\pi(\mu, \Sigma) \propto |\Sigma|^{-\left(\frac{\nu_0+p}{2}+1\right)} \exp\left[-\frac{1}{2} \text{tr}(\Lambda_0 \Sigma^{-1}) - \frac{k_0}{2} (\mu - \mu_0)' \Sigma^{-1} (\mu - \mu_0)\right]$$

The posterior density is  $\pi(\mu, \Sigma | \tilde{X}) =$  a Normal-Inverse Wishart density with parameters

$$\mu_n = \frac{k_0}{k_0+n} \mu_0 + \frac{n}{k_0+n} \bar{X},$$

$$k_n = k_0 + n, \quad \nu_n = \nu_0 + n,$$

$$\Lambda_n = \Lambda_0 + S + \frac{k_0 n}{k_0+n} (\bar{X} - \mu_0)(\bar{X} - \mu_0)'$$

where,  $S$  is the correlated sum of square matrix, i.e.,

$$S = \sum_{i=1}^n (x_i - \bar{X})(x_i - \bar{X})'$$

The marginal posterior distr. of  $\mu$  is a multivariate t-distr.

Regression Model:-  $y_i, x_{i1}, \dots, x_{ip}; i=1, 2, \dots, n.$

$$\tilde{y} = (y_1, \dots, y_n)', \quad X = ((x_{ij}))_{n \times p}$$

$$\tilde{y} = X\beta + \underline{\epsilon} \quad ; \text{ where, } \epsilon_1, \epsilon_2, \dots, \epsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$E(y_i | \beta, X) = \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

$$y | \beta, \sigma^2 \sim f(y | \beta, \sigma^2) = N_n(X\beta, \sigma^2 I)$$

Taking  
[ $y = \tilde{y}$ ]

Consider the standard non-informative prior  $\pi(\beta, \sigma^2)$ ; where,

$$\pi(\beta, \sigma^2) \propto \frac{1}{\sigma^2}; \quad f(y | \beta, \sigma^2): \text{likelihood.}$$

$$\pi(\beta, \sigma^2 | y) \propto f(y | \beta, \sigma^2) \pi(\beta, \sigma^2)$$

$$= \frac{1}{(\sigma^2)^{n/2+1}} \exp \left[ -\frac{1}{2\sigma^2} (y - X\beta)' (y - X\beta) \right]$$

$$= (\sigma^2)^{-n/2-1} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) + \underbrace{(y - X\hat{\beta})' (y - X\hat{\beta})}_{(n-p)\delta^2, \text{ say}} \right\} \right]$$

$$\text{where, } \hat{\beta} = (X'X)^{-1} X'y.$$

$$\pi(\beta | \sigma^2, y) \propto f(y | \beta, \sigma^2) \pi(\beta)$$

$$\propto \exp \left[ -\frac{1}{2\sigma^2} (\beta - \hat{\beta})' V^{-1} (\beta - \hat{\beta}) \right]; \text{ where } V = (X'X)^{-1}$$

$$\text{which is } N_p(\hat{\beta}, \sigma^2 V).$$

$$\pi(\sigma^2 | y) = \int \pi(\beta, \sigma^2 | y) d\beta$$

$$= \text{constant} \times (\sigma^2)^{-n/2-1} \exp \left[ -\frac{1}{2\sigma^2} (y - X\hat{\beta})' (y - X\hat{\beta}) \right]$$

$$\times |(X'X)^{-1} \sigma^2|^{1/2}$$

$$\propto (\sigma^2)^{-\frac{n-p}{2}-1} \exp \left[ -\frac{1}{2\sigma^2} (n-p)\delta^2 \right]$$

which is a scaled inverse  $\chi^2(n-p, \delta^2)$  distribution;

$$\text{where, } \delta^2 = \frac{1}{n-p} (y - X\hat{\beta})' (y - X\hat{\beta}).$$

Posterior density of  $y$  is given by:

$$\pi(\beta|y) = \int \pi(\beta, \sigma^2|y) d\sigma^2$$

$$= \text{constant} \times \int (\sigma^2)^{-n/2-1} \exp\left(-\frac{A}{2\sigma^2}\right) d\sigma^2$$

$$\text{where, } A = (n-p)s^2 + (\beta - \hat{\beta})(X'X)(\beta - \hat{\beta})$$

$$\propto A^{-n/2} \quad [\text{As obtained earlier in } N(\mu, \sigma^2) \text{ example}]$$

$$\propto \left[ 1 + \frac{(\beta - \hat{\beta})X'X(\beta - \hat{\beta})}{(n-p)s^2} \right]^{-n/2}$$

which is a multivariate  $t$ -distr. with  $(n-p)$  d.f.

Location vector:  $\hat{\beta}$ , scale matrix:  $s^2(X'X)^{-1}$ .

The Multinomial Model:- This model is used to describe data for which each observation is one of  $k$  possible outcomes. Sample space =  $\{1, 2, \dots, k\}$ , say.

$X_1, X_2, \dots, X_n$  iid with  $P(X_i = j) = p_j$ ;  $j = 1, 2, \dots, k$ .  
 $p_i \geq 0 \forall j$  and  $\sum_{j=1}^k p_j = 1$ .

Likelihood function:  $f(x_1, \dots, x_n | p) = p_1^{n_1} p_2^{n_2} \dots p_k^{n_k}$ ;  $n_j$ : # of  $X_i$ 's equal to  $j$ .

Conjugate prior:- The conjugate prior distr. is Dirichlet( $\alpha_1, \alpha_2, \dots, \alpha_k$ );  $\alpha_j > 0 \forall j$ .  
 [Multivariate generalisation of Beta]

$$\pi(p_1, \dots, p_k | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\alpha_1 + \alpha_2 + \dots + \alpha_k)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\dots\Gamma(\alpha_k)} p_1^{\alpha_1-1} p_2^{\alpha_2-1} \dots p_k^{\alpha_k-1}; p_j \geq 0 \forall j \text{ and } \sum_{j=1}^k p_j = 1.$$

Posterior distr. of  $p$  is Dirichlet( $\alpha_1 + n_1, \dots, \alpha_k + n_k$ ).

Note that  $\alpha_j = 1 \forall j$  corresponds to a uniform distribution.

Bayesian Nonparametrics:-  $P$  is a random probability distr.,  
 We need a prior prob. distr. of the "random" element  $P$ .

Challenges of Bayesian nonparametrics:-

Construction of prior  
 Posterior computation

Example (Bernoulli revisited):  $\mathcal{X} = \{0, 1\}$

$X_1, X_2, \dots, X_n$  iid  $\text{Bin}(1, p)$ .

$p^* = (1-p, p)$  where  $p \sim \text{Beta}(\alpha_0, \alpha_1)$

$$E(p^*) = (E(1-p), E(p)) = \left( \frac{\alpha_1}{\alpha_0 + \alpha_1}, \frac{\alpha_0}{\alpha_0 + \alpha_1} \right)$$

$$\alpha(\cdot) = \text{measure of } \mathcal{X} \text{ ; i.e., } \bar{\alpha} = \frac{\alpha(\cdot)}{\alpha(\mathcal{X})}$$

So, the prob. distr.  $\bar{\alpha}$  is the prior guess for random prob. distr.  $P$ .

Definition:- A random measure  $P$  on  $(\mathbb{R}, \mathcal{B})$  is said to possess a Dirichlet process distr.  $DP(\alpha)$  with base measure  $\alpha$ , if for every finite measurable partition  $A_1, \dots, A_k$  of  $\mathcal{X} (= \mathbb{R})$ ,  $(P(A_1), \dots, P(A_k)) \sim \text{Dirichlet}(\alpha(A_1), \dots, \alpha(A_k))$ .

Proposition:- If  $P \sim DP(\alpha)$ , then for any measurable set  $A$  and  $B$

$$E(P(A)) = \bar{\alpha}(A)$$

$$V(P(A)) = \frac{\bar{\alpha}(A) \bar{\alpha}(A^c)}{1 + \alpha(\mathcal{X})}$$

$$\text{Cov}(P(A), P(B)) = \frac{\bar{\alpha}(A \cap B) - \bar{\alpha}(A) \bar{\alpha}(B)}{1 + \alpha(\mathcal{X})}$$

Theorem:- The posterior distr. given a random sample from  $P$  where  $P \sim DP(\alpha)$  is  $DP\left(\alpha + \sum_{i=1}^n \delta_{x_i}\right)$ ; where  $\delta_c$  is a measure degenerate at  $c$ .

Note:- In real life problems, Dirichlet mixture distr. is highly useful to model the data.

## Consistency of Posterior Distribution:-

$$\tilde{X}_n = (X_1, \dots, X_n) \sim f(\tilde{x}_n | \theta), \theta \in \mathbb{H} \subset \mathbb{R}^p.$$

$\pi(\theta)$  is prior density.

$\pi(\theta | \tilde{X}_n)$  is posterior density.

$\pi(\cdot | \tilde{X}_n)$  is the corresponding posterior distr. .

• Definition:- The sequence of the posterior distr.  $\pi(\cdot | \tilde{X}_n)$  is said to be consistent at some  $\theta_0 \in \mathbb{H}$  if for every neighbourhood  $U$  of  $\theta_0$ ,  $\pi(U | \tilde{X}_n) \rightarrow 1$  as  $n \rightarrow \infty$  with prob. 1 w.r.t. the distr. under  $\theta_0$ .

• For a real parameter  $\theta$ , consistency at  $\theta_0$  can be proved by showing  $E(\theta | \tilde{X}_n) \rightarrow \theta_0$  and  $\text{Var}(\theta | \tilde{X}_n) \rightarrow 0$  w.p.1.

• Example:  $X_1, \dots, X_n$  be iid Bernoulli ( $\theta$ );  $\theta \sim \text{Beta}(\alpha, \beta)$  (prior)

Posterior distr. is  $\text{Beta}\left(\sum_{i=1}^n X_i + \alpha, n - \sum_{i=1}^n X_i + \beta\right)$  with

$$E(\theta | \tilde{X}_n) = \frac{\sum_{i=1}^n X_i + \alpha}{n + \alpha + \beta} \quad \text{and} \quad \text{Var}(\theta | \tilde{X}_n) = \frac{\left(\sum_{i=1}^n X_i + \alpha\right)\left(n - \sum_{i=1}^n X_i + \beta\right)}{(\alpha + \beta + n)^2 (\alpha + \beta + n + 1)}$$

As,  $\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \theta_0$  w.p.1 under  $\theta_0$  (SLLN)

It follows that  $E(\theta | \tilde{X}_n) \rightarrow \theta_0$  and  $\text{Var}(\theta | \tilde{X}_n) \rightarrow 0$  w.p.1 under  $\theta_0$ .

This implies that the posterior is consistent.

An important result related to consistency is the robustness of posterior inference w.r.t. choice of prior.

Let  $X_1, \dots, X_n$  be iid and  $\pi_1$  and  $\pi_2$  are two prior densities which are positive and continuous at  $\theta_0$ , an interior point of  $\mathbb{H}$  such that the corresponding posterior distr.  $\pi_1(\cdot | \tilde{X}_n)$  and  $\pi_2(\cdot | \tilde{X}_n)$  are both consistent at  $\theta_0$ ,

$$\int_{\mathbb{H}} |\pi_1(\theta | \tilde{X}_n) - \pi_2(\theta | \tilde{X}_n)| d\theta \rightarrow 0$$

$$\Leftrightarrow \sup_A |\pi_1(A | \tilde{X}_n) - \pi_2(A | \tilde{X}_n)| \rightarrow 0$$

On choice of Priors:- Given data, a Bayesian needs the likelihood  $p(x|\theta)$  and the prior  $\pi(\theta)$ . For many standard problems, the likelihood is known either from past experience or by conventions. To derive the Bayesian engine, one would need an appropriate prior.

Ideally one wants to choose a prior or a class of priors reflecting one's prior knowledge and belief about unknown parameters or about different hypotheses. This is a subjective choice. Choice of subjective priors, usually called elicitation of priors, is still rather difficult. We will discuss it in details later. But at this moment, let us see why it is difficult. Empirical studies have shown that experience and maturity helps the person in quantifying uncertainty about an event in the form of a probability. However assigning a fully specified prob. distr. to an unknown parameter is difficult even when the parameter has a physical meaning like length and, or, breadth of some article. In such cases, it may be realistic to expect elicitation of prior mean or variance to some other prior quantities but not a full specification of the distribution.

It is much more common to use what are called objective priors. When very little prior information is available, objective priors are also called non-informative prior. The term "non-informative prior" is no longer in favour among Bayesian because a complete lack of information is hard to define. There are standard algorithms to the construction of priors with low information. Objective priors are typically improper but have proper posterior distribution. They are suitable for estimation problems and also for testing problems where both null and alternative hypothesis have the same distributions. The objective prior need to be suitably modified for sharp null hypothesis.



## Different methods for construction of Objective Priors:-

We may do one of the following things to construct objective prior under general regularity condition:

1. Define a uniform distn. that takes into account the geometry of parameter space.
2. Minimize a suitable measure of information in the prior.
3. Choose a prior with same form of frequentist ideas because a prior with little information should lead to inference that is similar to frequentist inference.

- To fully define these methods, we have to specify the geometry in (1), the measure of information (2) and the frequentist ideas that are to be defined in (3). It is striking fact that (1) and (2) both lead to Jeffrey's prior, namely,  
$$\pi(\theta) = \sqrt{|I(\theta)|}$$
, where,  $I(\theta)$  is the Fisher Information Matrix. In the one dimensional case (3) leads to Jeffrey's prior.

- It may be noted that some statistical models possess additional structure. Some are exponential families, some are location-scale families or more generally families invariant under a group of transformation.

- For each of these special classes, there is a different choice of objective priors. The objective priors for the exponential families came from the class of conjugate priors. In the class of location-scale families with scale parameter  $\sigma$ ; the common objective prior is so called right invariant Haar measure,  $\pi_1(\mu, \sigma) = 1/\sigma$  and the Jeffrey's prior is the left invariant Haar measure  $\pi_2(\mu, \sigma) = \frac{1}{\sigma^2}$ . There are strong reasons for preferring  $\pi_1$  to  $\pi_2$ .

## The Uniform Prior and its criticism:-

The first objective prior ever to be used is the uniform distribution over a compact interval. A common argument based on "ignorance" seems to have been that if we know nothing about  $\theta$  why should I attach more density to one point than another. The principle of ignorance has been criticized by many. Essentially, the criticism is based on an invariance argument.

Let  $\eta = \psi(\theta)$  be a one-to-one function of  $\theta$ . If we know nothing about  $\theta$ , then we know nothing about  $\eta$  also. So the principle of ignorance applied to  $\eta$  will imply our prior  $\pi$  is uniform; just as it had led to a uniform prior for  $\theta$ . But this leads to contradiction.

To see this suppose  $\psi$  is differentiable and  $p(\eta) = c$  on  $\psi(\Theta)$  then the prior  $p^*(\theta)$  for  $\theta$  is  $p^*(\theta) = p(\eta) |\psi'(\theta)|$  i.e.,  $p^*(\theta) = c |\psi'(\theta)|$  which is not a constant in general.

This argument also leads to an invariance principle. Suppose we have an algorithm that produces non-informative priors for both  $\theta$  and  $\eta$ . Then these priors  $p^*(\theta)$  and  $p(\eta)$  should be connected by the equation;

$$p^*(\theta) = p(\eta) |\psi'(\theta)|$$

i.e., a non-informative prior should be invariant under one-to-one differentiable transformations.

Some arguments in favour of Jeffrey's prior as a non-informative

Prior:- Ghoshal, Ghosh, Ramamurthy (1997) constructed a uniform distn. taking into account the topology of the parameter space and showed that the uniform distn. is the Jeffrey's prior which automatically satisfy the invariance requirement.

In their construction one takes a compact subset of the parameter space and approximates this by a finite set of points in the Hellinger metric

$$d(P_\theta, P_{\theta'}) = \left[ \int (\sqrt{P_\theta} - \sqrt{P_{\theta'}})^2 dx \right]^{1/2}$$

where,  $P_\theta$  and  $P_{\theta'}$  are the densities of  $P_\theta$  and  $P_{\theta'}$ . One thus puts a discrete uniform distr. on the approximating finite set of points and let the degree of approximation tend to zero. then the corresponding discrete uniform converge to Jeffrey's prior.

$K$  is a compact metric space with a metric  $f$ .  $S \subset K$  is called  $\epsilon$ -dispersed if  $f(x, y) \geq \epsilon \quad \forall \quad x, y \in S, x \neq y$ .

A maximal  $\epsilon$ -dispersed set will be called an  $\epsilon$ -net.

An  $\epsilon$ -net with maximum possible cardinality is said to be an  $\epsilon$  lattice is called the packing number (or  $\epsilon$ -capacity) of  $K$  and is denoted by  $D(\epsilon, K)$ .

For  $X \subset K$  define  $P_\epsilon(X) = \frac{D(\epsilon, X)}{D(\epsilon, K)}$ .

In the paper, in place of  $f$  one takes the Hellinger metric then taking a compact subset  $K$  of  $(H)$  one has for  $\mathcal{Q} \subset K$

$$\lim_{\epsilon \rightarrow 0} \frac{D(\epsilon, \mathcal{Q})}{D(\epsilon, K)} = \frac{\int_{\mathcal{Q}} \sqrt{\det I(\theta)} d\theta}{\int_K \sqrt{\det I(\theta)} d\theta}$$

## Jeffrey's Prior as a

## Probability matching prior :-

One should expect an objective prior with low information to provide inference similar to that based on uniform prior for  $\theta$  in  $N(\theta, 1)$ , where the Bayesian and frequentist answers were exactly identical. It may be recalled that in this problem, the posterior distr. of  $\theta - \bar{X}$  given  $\underline{X}$  is identical with the frequentist distr. of  $\theta - \bar{X}$  given  $\theta$ .

In general case, we will not get exactly the same distr. but only up to  $O(1/n)$ . We give a precise definition of probability matching prior for a single parameter below.

Let  $X_1, X_2, \dots, X_n$  be i.i.d.  $p(x|\theta)$ ,  $\theta \in \Theta \subset \mathbb{R}$ . Assume regularity condition of the parameter with normal  $(\hat{\theta}, (nI(\hat{\theta}))^{-1})$  as the leading term. For  $0 < \alpha < 1$ , choose  $\theta_\alpha(x)$  depending on the prior  $p(\theta)$  s.t.

$$P(\theta \leq \theta_\alpha(x) | X) = 1 - \alpha + O_p(1/n).$$

We say  $p(\theta)$  is probability matching (to the 1st order) if

$$P_\theta(\theta \leq \theta_\alpha(x)) = 1 - \alpha + O(1/n)$$

Uniformly on compact set of  $\theta$  using normal approximation,

$$P(\theta < \hat{\theta} + \frac{z_\alpha}{\sqrt{n}} | X) = P_\theta(\theta < \hat{\theta} + \frac{z_\alpha}{\sqrt{n}}) + O_p(1/n).$$

under  $\theta$  (uniformly on compact subsets of  $\Theta$ ).

Probability matching priors :- Ref: Dutta & Mukherjee (2009).

In this approach, a common goal is to find matching prior for posterior quantities, we want

$$P(\theta < \theta^{(1-\alpha)}(\pi, X)) = 1 - \alpha + O(n^{-1/2})$$

for each  $\alpha$ ;  $0 < \alpha < 1$ , where  $\theta^{(1-\alpha)}(\pi, X)$  is the  $(1-\alpha)$ th posterior quantile of  $\theta$  under priors  $\pi(\cdot)$  given data  $X$ .

In order to do this first one finds an expression for the posterior density of  $\sqrt{n}(\theta - \hat{\theta})$ .

$$\pi(\theta|x) = \frac{\pi(\theta) \exp(nl(\theta))}{\int \pi(\theta) \exp(nl(\theta)) d\theta}$$

$$= \frac{\pi(\theta) \exp(n(l(\theta) - l(\hat{\theta})))}{\int \pi(\theta) \exp(n(l(\theta) - l(\hat{\theta}))) d\theta}$$

Here  $l(\theta) = \frac{1}{n} \sum_{i=1}^n \log(f(x_i, \theta))$

let  $h = \sqrt{n}(\theta - \hat{\theta})$

$$\pi^*(h|x) = \frac{b(h, x)}{\int b(h, x) dh}; \text{ where,}$$

$$b(h, x) = \pi\left(\hat{\theta} + \frac{h}{\sqrt{n}}\right) \exp\left\{n\left(l\left(\hat{\theta} + \frac{h}{\sqrt{n}}\right) - l(\hat{\theta})\right)\right\}$$

Then do Taylor expansion about  $\hat{\theta}$  in both these terms:

this will finally give

$$\pi^*(h|x) \dots \phi_h(0, c^{-1}) \left[1 + \frac{1}{\sqrt{n}} g_1(h, x) + \frac{1}{n} g_2(h, x)\right] + o(1/n);$$

where,  $\phi_h(0, c^{-1})$  is the density of a normal distr. with mean zero and variance as the inverse of per observation observed Fisher Information matrix ( $c$ ), evaluated at  $h$ .

Then in the next step one finds an expression for  $\theta^{(1-\alpha)}(\pi, x)$ , the  $(1-\alpha)$ th posterior quantile of  $\theta$  under prior  $\pi(\cdot)$ . This can be obtained by using the expansion of the posterior already estimated. In general, it will be of the form:

$$\theta^{(1-\alpha)}(\pi, x) = \hat{\theta} + \frac{1}{\sqrt{n}} c(x) + \left[ z_\alpha + \frac{1}{\sqrt{n}} \left( \quad \right) + \frac{1}{n} \left( \quad \right) \right]$$

The last and very crucial step is to calculate an expression for the frequentist coverage probability  $P_\theta(\theta \leq \theta_1^{(1-\alpha)}(\pi, x))$

$$P(\theta \leq \theta_1^{(1-\alpha)}(\pi, x) | x) = 1 - \alpha + o(1/n)$$

Using the shrinkage argument, directed by J.K. Ghosh one can find an expression for this quantity. For details, see Dutta and Mukherjee (2004). Using this technique, one ultimately gets —

$$P_{\theta}(\theta \leq \theta^{(1-\alpha)}(\pi, x)) = (1-\alpha) + \frac{1}{\sqrt{n}} \cdot \frac{\phi(z_{\alpha})}{\pi(\theta)} A_1(\pi, \theta) + \frac{1}{n} \cdot \frac{z_{\alpha} \phi(z_{\alpha})}{\pi(\theta)} A_2(\pi, \theta) + O(1/n),$$

for specific function  $A_1$  and  $A_2$ .

In this case,  $A_1(\pi, \theta) = \frac{d}{d\theta} \left\{ \frac{\pi(\theta)}{I^{1/2}} \right\}$

So, for a prior to be first order probability matching, one needs to have  $A_1(\pi, \theta) = 0$ ; i.e.,  $\frac{d}{d\theta} \left\{ \frac{\pi(\theta)}{I^{1/2}} \right\} = 0$ .

The solution of this equation comes out as

$$\pi(\theta) \propto \sqrt{I(\theta)}, \text{ the well known Jeffreys' prior.}$$

Jeffrey's prior as a minimizer of information:-

The amount of information to be expected from an experiment about some quantities of interest naturally depends on the available prior knowledge; the more prior information available, the less information can be learned from the data. An infinitely large experiment would eventually provide all missing information, thus it is possible to obtain a measure of the amount of missing information as a limiting form of a function of the prior distr. Using an idea of Lindley, Bernardo suggested that a Kullback-Leibler divergence between prior and posterior namely

$J(P(\theta), X) = E \left( \log \frac{P(\theta|X)}{P(\theta)} \right)$  as the amount of information which may be expected to be provided by an experiment about the value of  $\theta$ . It is easy to see that when the prior  $P(\theta)$  is nearly degenerate at a point, so will be the posterior, then  $J$  will be near 0.

On the other hand, if  $P(\theta)$  is rather diffuse  $P(\theta|X)$  will differ a lot from  $P(\theta)$  at least for large or moderate  $n$  because  $P(\theta|X)$  would be approximately normal with mean  $\hat{\theta}$  and variance of the order  $O(1/n)$ . The substantial difference in priors and posterior will be reflected in a large value for  $J$ . To sum up:  $J$  is small when  $P(\theta)$  is nearly degenerate and large when  $P$  is diffuse, i.e.,  $J$  combines how diffuse is the priors. Therefore it makes sense to maximize  $J$  w.r.t. the priors  $P(\theta)$ .

Bernardo suggested one should not work with sample size  $n$  of the given data and maximize  $J$  for this  $n$ . For one thing, this would be technically forbidding in most cases and more importantly, the functional  $J$  is expected to be a nice function of the priors only asymptotically.

We will do the maximisation asymptotically. One may approximate  $J$  by  $\hat{J}$  by using idea of posterior normality, i.e., by replacing  $p(\theta|x)$  by  $\hat{p}(\theta|x)$  in the definition of  $J$ ,  $\hat{p}(\theta|x)$  being the density of the approximately normal dist.

With this approximation, one gets

$$\hat{J} = \left\{ -\frac{d}{2} \log(2\pi) - \frac{d}{2} + \frac{d}{2} \log n \right\} + \int \log(\sqrt{I(\theta)} P_i(\theta)) d\theta - \int_{K_i} \log(P_i(\theta)) P_i(\theta) d\theta + O_p(1).$$

$(-J)$  is a measure of information in a priors.

To maximize  $J$  w.r.t.  $p(\theta)$ ; we will try asymptotic maximization of  $J$ .

Towards that; let us concentrate on a compact subset  $K_i$  of  $(H)$ , you can show that

$$J(p(\theta), X) = \int_{(H)} \left\{ \int_x \left[ \int_{(H')} \frac{\log(p(\theta'|x))}{p(\theta')} p(\theta'|x) d\theta' \right] p(x|\theta) dx \right\} p(\theta) d\theta$$

We will first assume that

$$J(p(\theta), X) \approx \hat{J}(\hat{p}(\theta), X)$$

We will assume that  $J$  can be approximated well approximately by replacing  $p(\theta'|x)$  by  $\hat{p}(\theta'|x)$ ; where  $\hat{p}(\theta'|x)$  is the density of a  $N(\hat{\theta}, \frac{1}{nI(\hat{\theta})})$  distn. evaluated at  $\theta'$ .

Restricted to compact set  $K_i$

$$\hat{J}(\hat{p}(\theta), X) = \int_{K_i} \left\{ \int_x \left[ \int_{\mathbb{R}} \log \frac{\hat{p}_i(\theta'|x)}{p_i(\theta')} \hat{p}_i(\theta'|x) d\theta' \right] p(x|\theta) dx \right\} p_i(\theta) d\theta$$

$$\theta' \approx N\left(\hat{\theta}, \frac{1}{nI(\hat{\theta})}\right)$$

$$\hat{p}_i(\theta'|x) = \frac{\sqrt{nI(\hat{\theta})}}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} nI(\hat{\theta})(\theta' - \hat{\theta})^2\right)$$

$$\log \hat{p}_i(\theta'|x) = \frac{1}{2} \log n + \frac{1}{2} \log I(\hat{\theta}) - \frac{1}{2} \log 2\pi - \frac{1}{2} nI(\hat{\theta})(\theta' - \hat{\theta})^2$$

$\hat{J}(\hat{p}(\theta), X)$  reduces to

$$\approx \int_{K_i} \left[ \int_x -\log p_i(\hat{\theta}) p(x|\theta) dx \right] p_i(\theta) d\theta$$

$$\approx - \int_{K_i} \log(p_i(\hat{\theta})) p_i(\theta) d\theta$$

$$\begin{aligned} \text{Big integral} &\approx \left(-\frac{1}{2} \log 2\pi - \frac{1}{2} + \frac{1}{2} \log n\right) + \int_{K_i} \log(\sqrt{I(\hat{\theta})}) p_i(\theta) d\theta \\ &\quad - \int_{K_i} \log(p_i(\hat{\theta})) p_i(\theta) d\theta \\ &= c + \int_{K_i} \log\left(\frac{\sqrt{I(\hat{\theta})}}{p_i(\hat{\theta})}\right) p_i(\theta) d\theta \end{aligned}$$



The prior maximizes this quantity if

$$p_i(\theta) = \begin{cases} c_i \sqrt{I(\theta)} & \text{if } \theta \in K_i \\ 0 & \text{on } \text{on} \end{cases}$$

where,  $c_i$  is a normalising constant s.t.  $c_i \sqrt{I(\theta)}$  is a prob. density on  $K_i$ .

Bayesian asymptotics: (Posterior consistency and Posterior normality)

The Bernstein-Van Mises theorem about posterior on normality says the following:

Suppose  $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} f(x|\theta); \theta \in \Theta \subset \mathbb{R}$

where,  $\Theta$  is an open subset of  $\mathbb{R}$ .

Let  $\theta_0 \in \Theta$ ; where  $\Theta$  is an open subset of  $\mathbb{R}$ .

Let  $\hat{\theta}_n$  be a strongly consistent solution of the likelihood equation. Assume all the usual regularity conditions needed for proving asymptotic normality of consistent roots of MLE.

(i) for any  $\delta > 0$

$P_{\theta_0} \left[ \text{for some } \epsilon > 0, \sup_{|\theta - \theta_0| > \epsilon} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)] < -\epsilon, \text{ for all sufficiently large } n \right] = 1.$

(ii) The prior has a density  $\pi(\theta)$  w.r.t. Lebesgue measure, which is continuous and positive at  $\theta_0$ .

Let  $\pi_n^*(t | X_1, X_2, \dots, X_n)$  be the posterior density of  $t = \sqrt{n}(\theta - \hat{\theta}_n)$

Then 
$$\lim_{n \rightarrow \infty} \int \left| \pi_n^*(t | X_1, \dots, X_n) - \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2} t^2 I(\theta_0)} \right| dt$$

$= 0$  with probability 1, under  $P_{\theta_0}$ .

Bernstein-Van-Mises Theorem: - A standard result in the asymptotic theory of MLE on its asymptotic normality. Its Bayesian parallel is a result often referred to as Bernstein-Van-Mises theorem - a result on the asymptotic normality of the posterior distr. We consider the case of a real parameter  $\theta$ ,  $\theta \in \mathbb{H} \subset \mathbb{R}$ , where  $\mathbb{H}$  is open.

For simplicity, let us take  $\mathbb{H} = \mathbb{R}$ . But the result goes through exactly as stated below even though  $\mathbb{H}$  is an open interval in  $\mathbb{R}$ .

Suppose  $X_1, \dots, X_n \stackrel{iid}{\sim} P_\theta$ ,  $\theta \in \mathbb{R}$ .

$P_\theta$  has the density  $f(x|\theta)$ . We make the following regularity assumptions. Fix  $\theta_0 \in \mathbb{H}$  (regarded as the true value of the parameter)

Notation:  $l(\theta, x) = \log f(x|\theta)$   
 $L_n(\theta) = \log \prod_{i=1}^n f(x_i|\theta) = \sum_{i=1}^n l(\theta, x_i)$ .

The assumptions are as follows:

(A1) The set  $\{x: f(x|\theta) > 0\}$  is the same for all  $\theta \in \mathbb{R}$ .

(A2)  $l(\theta, x) = \log f(x|\theta)$  is thrice differentiable w.r.t.  $\theta$  in the neighbourhood  $(\theta_0 - \delta, \theta_0 + \delta)$ . If  $l'$ ,  $l''$  and  $l'''$  be the 1st, 2nd and 3rd order derivatives then  $E_{\theta_0}(l'(\theta_0))$  and  $E_{\theta_0}(l''(\theta_0))$  are both finite, and

Sup  $|l'''(\theta, x)| \leq M(x)$  and  $E_{\theta_0}(M(x_1)) < \infty$ .

$\theta \in (\theta_0 - \delta, \theta_0 + \delta)$

(A3) Interchange of the order of expectation w.r.t.  $P_{\theta_0}$  and differentiation at  $\theta_0$  are justified, so that

$E_{\theta_0}(l'(\theta_0)) = 0$ ,  $E_{\theta_0}(l''(\theta_0)) = -E_{\theta_0}(l'(\theta_0))^2 = -I(\theta_0)$  (say)

(A4)  $0 < I(\theta_0) < \infty$ .

(A5) For any  $\delta > 0$

$$P_{\theta_0} \left[ \text{for some } \epsilon > 0, \sup_{|\theta - \theta_0| > \delta} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)] \leq -\epsilon \text{ for all sufficiently large } n \right] = 1.$$

(A6)  $\theta \sim \pi(\theta)$  where,  $\pi(\theta)$  is continuous and positive at  $\theta_0$ .  
 Let  $\hat{\theta}_n$  be a strongly consistent sequence of soln. of likelihood equation.

Let  $\pi_n^*(t | x_1, \dots, x_n)$  denote the posterior density of  $t = \sqrt{n}(\theta - \hat{\theta}_n)$ .

Then,

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} \left| \underbrace{\pi_n^*(t | x_1, \dots, x_n)}_{\substack{\text{random quantity} \\ \text{that's why prob.} \\ \text{comes into play}}} - \underbrace{\frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2 I(\theta_0)}}_{\substack{\text{density of } N(0, \frac{1}{I(\theta_0)}) \\ \text{with prob. 1 under} \\ P_{\theta_0}}} \right| dt = 0$$

The evaluation is done considering  $\theta_0$  to be the true value.

## Asymptotic Normality of Posterior distribution:-

Let  $\tilde{\theta}_n$  be the posterior mode.  $\underline{x}_n = (x_1, \dots, x_n)$

$\pi(\theta | \underline{x}_n)$ : posterior density

Under suitable regularity conditions, a Taylor series expansion of  $\log \pi(\theta | \underline{x}_n)$  at  $\hat{\theta}_n$  gives —

$$\log \pi(\theta | \underline{x}_n) = \underbrace{\log \pi(\tilde{\theta}_n | \underline{x}_n)}_{\text{free of } \theta} + (\theta - \tilde{\theta}_n) \frac{\partial}{\partial \theta} \log \pi(\theta | \underline{x}_n) \Big|_{\tilde{\theta}_n} - \frac{1}{2} (\theta - \tilde{\theta}_n)' \tilde{I}_n (\theta - \tilde{\theta}_n) + \dots$$

where,  $\tilde{I}_n$  is a  $p \times p$  matrix defined as follows:

$$\log \pi(\theta | \underline{x}_n) \approx \log \pi(\hat{\theta}_n | \underline{x}_n) - \frac{1}{2} (\theta - \tilde{\theta}_n)' \tilde{I}_n (\theta - \tilde{\theta}_n)$$

$$\pi(\theta | \underline{x}_n) \propto \exp \left[ -\frac{1}{2} (\theta - \tilde{\theta}_n)' \tilde{I}_n (\theta - \tilde{\theta}_n) \right] \quad (\text{approximately})$$

a  $N_p(\tilde{\theta}_n, \tilde{I}_n^{-1})$  density for  $\theta$ .

$$\tilde{I}_n = \left( -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\underline{x}_n | \theta) \right)_{\tilde{\theta}_n = \text{MLE}}$$

Fisher Information matrix (generalised observed Fisher Information matrix). is the observed

Result: Under suitable regularity assumptions, for large  $n$ , the posterior distr. of  $\theta$  can be approximated by any one of the normal distr.  $N_p(\tilde{\theta}_n, \tilde{I}_n^{-1})$  or  $N_p(\tilde{\theta}_n, I^{-1}(\tilde{\theta}_n))$ . In particular, under suitable regularity conditions, the posterior distr. of  $\tilde{I}_n^{1/2} (\theta - \tilde{\theta}_n)$  gives  $\underline{x}_n$  converges to  $N_p(0, I)$ ; where,  $I =$  Identity matrix of order  $p$ .

— We now formally state a theorem giving a set of regularity conditions under which asymptotic normality of posterior distr. holds.

— Let  $x_1, \dots, x_n$  be iid obsn with common distr.  $P_\theta$  possessing a density  $f(x|\theta)$  where  $\theta \in \mathbb{H}$ , an open subset of  $\mathbb{R}$ . We fix some  $\theta_0 \in \mathbb{H}$  which may be regarded as the "true" value of the parameter as the probability statements are all made under  $\theta_0$ .

Let  $l(\theta, x) = \log f(x|\theta)$

$L_n(\theta) = \log \prod_{i=1}^n f(x_i|\theta) = \sum_{i=1}^n l(\theta, x_i)$  : log-likelihood function.

For a function,  $h(\cdot)$ , let  $h^{(i)}$  denotes the  $i^{th}$  derivative of  $h$ .

This is also known as Bernstein-Von-Mises theorem or Bayesian Central Limit theorem.

The assumptions are as follows:

(A1) The set  $\{x : f(x|\theta) > 0\}$  is the same for all  $\theta \in \Theta$

(A2)  $l(\theta, x)$  is thrice differentiable w.r.t.  $\theta$  in a neighbourhood  $(\theta_0 - \delta, \theta_0 + \delta)$  of  $\theta_0$ . the expectations  $E_{\theta_0} l^{(1)}(\theta_0, x_1)$  and  $E_{\theta_0} l^{(2)}(\theta_0, x_1)$  are both finite and

Sup  $|l^{(3)}(\theta, x)| \leq M(x)$  and  $E_{\theta_0}(M(x_1)) < \infty$ .

$\theta \in (\theta_0 - \delta, \theta_0 + \delta)$

(A3) Interchange of the order of expectation w.r.t.  $P_{\theta_0}$  and differentiation at  $\theta_0$  are justified so that

$E_{\theta_0} l^{(1)}(\theta_0, x_1) = 0$ ,

$E_{\theta_0} l^{(2)}(\theta_0, x_1) = -E_{\theta_0} [l^{(1)}(\theta_0, x_1)]^2 = -I(\theta_0)$ ,

where,  $I(\theta_0)$  is the Fisher Information number per unit observation is positive and finite.

(A4) For any  $\delta > 0$  with  $P_{\theta_0}$  - probability one

Sup  $\frac{1}{n} [\ln(\theta) - \ln(\theta_0)] < -\epsilon$

for some  $\epsilon > 0$  and for all sufficiently large  $n$ .

Remark:- We assume that there is a strongly consistent solution  $\hat{\theta}_n$  to the likelihood equation  $L_n^{(1)}(\theta) = 0$ , i.e.,  $\exists$  a sequence of statistics  $\hat{\theta}_n \ni$  with  $P_{\theta_0}$  - probability one  $\hat{\theta}_n$  satisfies the likelihood equation for sufficiently large  $n$  and  $\hat{\theta}_n \rightarrow \theta_0$ .

Theorem:- Suppose assumption (A1)-(A4) hold and  $\hat{\theta}_n$  is a strongly consistent solution of the likelihood equation, then for any prior density  $\pi(\theta)$  which is continuous and positive at  $\theta_0$ ,

$$\lim_{n \rightarrow \infty} \int_{\mathbb{R}} \left| \pi_n^*(t | x_1, \dots, x_n) - \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt = 0 \quad \text{--- (*)}$$

with  $P_\theta$  - probability one, where,  $\pi_n^*(t | x_1, \dots, x_n)$  is the posterior density of  $t = \sqrt{n}(\theta - \hat{\theta}_n)$  given  $x_1, \dots, x_n$ .

Also, under the same assumptions, (\*) holds with  $I(\theta_0)$  replaced by  $\hat{I}_n = -\frac{1}{n} L_n^{(2)}(\hat{\theta}_n)$ .

Proof:

Posterior density of  $\theta$  given  $x_1, \dots, x_n$  is  $\pi_n(\theta | x_1, \dots, x_n)$  which equals to

$$\pi_n(\theta | x_1, \dots, x_n) = \frac{\prod_{i=1}^n f(x_i | \theta) \pi(\theta)}{\int \prod_{i=1}^n f(x_i | \theta') \pi(\theta') d\theta'}$$

[ Putting  $\theta = \hat{\theta}_n + \frac{t}{\sqrt{n}}$  ]

Posterior distr. of  $t = \sqrt{n}(\theta - \hat{\theta}_n)$  is

$$\pi_n^*(t | x_1, \dots, x_n) = \frac{\pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \prod_{i=1}^n f\left(x_i | \hat{\theta}_n + \frac{t}{\sqrt{n}}\right)}{\int \prod_{i=1}^n f\left(x_i | \hat{\theta}_n + \frac{t'}{\sqrt{n}}\right) \pi\left(\hat{\theta}_n + \frac{t'}{\sqrt{n}}\right) dt'}$$

$$\therefore \pi_n^*(t | x_n) = \frac{\pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \left[ \prod_{i=1}^n f\left(x_i + \hat{\theta}_n + \frac{t}{\sqrt{n}}\right) / \prod_{i=1}^n f\left(x_i | \hat{\theta}_n\right) \right]}{\int \pi\left(\hat{\theta}_n + \frac{t'}{\sqrt{n}}\right) \prod_{i=1}^n \frac{f\left(x_i + \hat{\theta}_n + \frac{t'}{\sqrt{n}}\right)}{f\left(x_i | \hat{\theta}_n\right)} dt'}$$

$$= \frac{\pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left[ L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n) \right]}{\int \pi\left(\hat{\theta}_n + \frac{t'}{\sqrt{n}}\right) \exp\left\{ L_n\left(\hat{\theta}_n + \frac{t'}{\sqrt{n}}\right) - L_n(\hat{\theta}_n) \right\} dt'}$$

$$= C_n^{-1} \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{ L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n) \right\}$$

[  $L_n(\theta) = \log \pi f(x_i | \theta)$   
 $= \sum \log f(x_i | \theta)$  ]

Need to show,

$$\textcircled{1} \dots \int_{\mathbb{R}} \left| c_n^{-1} \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n)\right\} - \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt \rightarrow 0 \text{ a.s. } [P_{\theta_0}]$$

(almost surely)

We first note that it is enough to show:

$$\textcircled{2} \dots \int_{\mathbb{R}} \left| \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n)\right\} - \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt \rightarrow 0 \text{ a.s.}$$

(under a bracket)  $g_n(t)$ , say

To see this, note that  $\textcircled{1}$  equals to

$$\textcircled{3} \dots c_n^{-1} \left[ \int_{\mathbb{R}} \left| \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n)\right\} - c_n \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt \right]$$

and as  $\textcircled{2}$  implies  $c_n \rightarrow \pi(\theta_0) \sqrt{\frac{2\pi}{I(\theta_0)}}$ , it is enough to show that the integral  $\int$  in  $\textcircled{3}$  goes to 0 a.s.

But  $0 \leq I \leq I_1 + I_2$  ; where,

$$I_1 = \int_{\mathbb{R}} \left| \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n)\right\} - \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt$$

$$I_2 = \int_{\mathbb{R}} \left| \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} - c_n \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2 I(\theta_0)} \right| dt$$

Now  $I_1 \rightarrow 0$  a.s. by  $\textcircled{2}$  and

$$I_2 = \left| \pi(\theta_0) - c_n \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} \right| \int_{\mathbb{R}} e^{-\frac{1}{2}t^2 I(\theta_0)} dt \rightarrow 0$$

as  $c_n \rightarrow \pi(\theta_0) \sqrt{\frac{2\pi}{I(\theta_0)}}$ .

Now we are to prove (2).

$$\int_{\mathbb{R}} |g_n(t)| dt \rightarrow 0 \text{ a.s. } [P_{\theta_0}]$$

where  $g_n(t) = \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left[\ln\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - \ln(\hat{\theta}_n)\right] - \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)}$ .

Note:  $\int_{\mathbb{R}} |g_n(t)| dt = \int_{A_1} |g_n(t)| dt + \int_{A_2} |g_n(t)| dt$ , where

$$A_1 = \left\{t : |t| < \delta_0 \sqrt{n}\right\} \text{ and } A_2 = \left\{t : |t| > \delta_0 \sqrt{n}\right\}$$

$\delta_0$  is approximately chosen small (+)ve number.

[We choose  $\delta_0$  so small that  $2\delta_0 < \delta$ ,  $|\theta - \theta_0| < 2\delta_0$ ,  $\pi(\theta) < 2\pi(\theta_0)$ ,  $\frac{\delta_0}{\theta} E_{\theta_0} M(x_1) < \frac{\pi(\theta_0)}{2}$ ]

Result 1:  $\int_{A_1} |g_n(t)| dt \rightarrow 0$  a.s.  $[P_{\theta_0}]$ ; where,  $A_1 = \left\{t : |t| < \delta_0 \sqrt{n}\right\}$

Note that with  $P_{\theta}$  - probability one, the following holds:

$$\left\{ \hat{\theta}_n \rightarrow \theta_0, -\frac{1}{n} \ln''(\hat{\theta}_n) \rightarrow I(\theta_0), \frac{1}{n} \sum_{i=1}^n M(x_i) \rightarrow E_{\theta_0}(M(x_1)) \right\} \text{ (*)}$$

We fix a sample sequence  $(x_1, x_2, \dots)$  for which (\*) hold and prove our result for this sequence

$$\ln\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - \ln(\hat{\theta}_n) = \frac{t^2}{2n} \ln''(\hat{\theta}_n) + \frac{1}{6} \left(\frac{t}{\sqrt{n}}\right)^2 \ln^{(3)}(\theta_n')$$

$$= -\frac{1}{2} t^2 \hat{I}_n + R_n(t), \text{ say, } \theta_n' \text{ lying between } \hat{\theta}_n + \frac{t}{\sqrt{n}} \text{ and } \hat{\theta}_n.$$

$$\text{and } \hat{I}_n = -\frac{1}{n} \ln''(\hat{\theta}_n).$$

$$\text{So, } \hat{I}_n = -\frac{1}{n} \ln''(\hat{\theta}_n) = \underbrace{-\frac{1}{n} \ln''(\theta_0)}_{\rightarrow I(\theta_0)} - \underbrace{\frac{1}{n} (\hat{\theta}_n - \theta_0) \ln^{(3)}(\theta_n'')}_{\rightarrow 0}$$

where,  $\theta_n''$  lying between  $\theta_0$  and  $\hat{\theta}_n$ .



$$\left[ \text{As } \left| \frac{1}{n} L_n^{(3)}(\theta_n) \right| \leq \frac{1}{n} \sum_{i=1}^n M(x_i) \neq n \text{ (for sufficiently large } n) \right. \\ \left. \rightarrow E_{\theta} (M(x_i)) < \infty \text{ by (A3)} \right]$$

$$\therefore \hat{I}_n \rightarrow I(\theta_0)$$

For each fixed  $t$ ,

$$\begin{aligned} |R_n(t)| &= \frac{|t|^3}{6\sqrt{n}} \left| \frac{1}{n} L_n^{(3)}(\theta_n) \right| \\ &\leq \frac{|t|^3}{6\sqrt{n}} \times \frac{1}{n} \sum_{i=1}^n M(x_i) \neq n \\ &\rightarrow 0 \end{aligned}$$

Thus for each fixed  $t$ ,  $g_n(t) \rightarrow 0$ .

For fix  $t \in A$ ,

$$\begin{aligned} |R_n(t)| &= \left( \frac{|t|}{\sqrt{n}} \right)^3 \times \frac{1}{6} \left| L_n^{(3)}(\theta_n) \right| \\ &\leq \delta_0 \cdot \frac{t^2}{6} \left| \frac{1}{n} L_n^{(3)}(\theta_n) \right| \\ &\leq \delta_0 \times \frac{t^2}{6} \times \frac{1}{n} \sum_{i=1}^n M(x_i) \neq n. \end{aligned}$$

$$\left[ \left| \hat{\theta}_n + \frac{t}{\sqrt{n}} - \theta_0 \right| < 2\delta_0, \quad \left| \hat{\theta}_n - \theta_0 \right| < \delta_0 \neq n \right]$$

$$\Rightarrow \left| \theta_n' - \theta_0 \right| < 2\delta_0 < \delta \neq n$$

$$\text{Now, } \frac{\delta_0}{6} \times \frac{1}{n} \sum_{i=1}^n M(x_i) < \frac{1}{4} \hat{I}_n \neq n$$

$$\left[ \begin{array}{l} \text{since LHS} \rightarrow \frac{\delta_0}{6} E_{\theta_0} M(x_i) \\ < \frac{1}{8} I(\theta_0) \end{array} \right]$$

$$\text{For } t \in A, \quad |R_n(t)| < \frac{1}{4} t^2 \hat{I}_n \neq n$$

$$\left[ \begin{array}{l} \text{RHS} \rightarrow \frac{1}{4} I(\theta_0) \end{array} \right]$$

$$\therefore L_n \left( \hat{\theta}_n + \frac{t}{\sqrt{n}} \right) - L_n(\hat{\theta}_n) \leq \frac{1}{2} t^2 \hat{I}_n + \frac{1}{4} t^2 \hat{I}_n = -\frac{1}{4} t^2 \hat{I}_n$$

$$\exp \left( \ln \left( \hat{\theta}_n + \frac{t}{\sqrt{n}} \right) - \ln(\hat{\theta}_n) \right) < e^{-\frac{1}{4} t^2 \hat{I}_n} < e^{-\frac{1}{8} t^2 I(\theta_0)}$$

$$\left[ \because \hat{I}_n > \frac{1}{2} I(\theta_0) \neq n \right]$$

$$\text{Then } |g_n(t)| \leq 2\pi(\theta_0) e^{-\frac{1}{8} t^2 I(\theta_0)} + \pi(\theta_0) e^{-\frac{1}{2} t^2 I(\theta_0)}$$

$$\text{We have proved that } \int_{A_1} |g_n(t)| dt \rightarrow 0 \text{ a.s. } [P_{\theta_0}]$$

Now we will prove that  $\int_{A_2} |g_n(t)| dt \rightarrow 0$  a.s.  $[P_{\theta_0}]$ .

Note,  $\int_{A_2} |g_n(t)| dt \leq \int_{A_2} \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) \exp\left\{L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n)\right\} dt$

$\leftarrow e^{-\frac{n\epsilon}{2}} + \int_{A_2} \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} dt$   
 $\rightarrow 0$  (normal tail) (3)

For  $t \in A_2$ ,

$$\frac{1}{n} \left\{ L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n) \right\} = \frac{1}{n} \left\{ L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\theta_0) \right\} + \frac{1}{n} \left\{ L_n(\theta_0) - L_n(\hat{\theta}_n) \right\}$$

$\left[ \text{For } t \in A_2, \left| \hat{\theta}_n + \frac{t}{\sqrt{n}} - \theta_0 \right| > \frac{\delta_0}{2} \forall n \right]$

$$\leq \sup_{|\theta - \theta_0| > \frac{\delta_0}{2}} \frac{1}{n} \left[ L_n(\theta) - L_n(\theta_0) \right] + \frac{1}{n} \left\{ \frac{1}{2} (\theta_0 - \hat{\theta}_n)^2 L_n''(\theta_n) + \frac{(\theta_0 - \hat{\theta}_n)^3}{6} L_n^{(3)}(\theta_n) \right\}$$

where,  $\theta_n$  lies between  $\theta_0$  and  $\hat{\theta}_n$ .

$< -\frac{\epsilon}{2} \forall n$

By assumption (A4), with  $P_{\theta_0}$ -probability one,

$\sup_{|\theta - \theta_0| > \frac{\delta_0}{2}} \frac{1}{n} \left[ L_n(\theta) - L_n(\theta_0) \right] < -\epsilon$  for some  $\epsilon > 0 \forall n$  (\*\*) (\*\*) in

Suppose the sample sequence  $(X_1, X_2, \dots)$  satisfies (\*\*) in addition to (\*),  
 Then  $\forall t \in A_2, \frac{1}{n} \left\{ L_n\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) - L_n(\hat{\theta}_n) \right\} < -\epsilon/2 \forall n$ .

Now from (3),

$$\int_{A_2} |g_n(t)| dt \leq \int_{A_2} \pi\left(\hat{\theta}_n + \frac{t}{\sqrt{n}}\right) e^{-n\epsilon/2} dt + \int_{A_2} \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} dt$$

$$\leq e^{-n\epsilon/2} \underbrace{\int_{\mathbb{R}} \pi(\theta) d\theta}_{\rightarrow 0} + \underbrace{\int_{A_2} \pi(\theta_0) e^{-\frac{1}{2}t^2 I(\theta_0)} dt}_{\rightarrow 0 \text{ (normal tail)}}$$

### Proof of 2nd part:-

$$\begin{aligned}
 \text{To show } & \int_{\mathbb{R}} \left| \pi^*(t | x_1, \dots, x_n) - \frac{\sqrt{\hat{I}_n}}{\sqrt{2\pi}} e^{-\frac{1}{2} t^2 \hat{I}_n} \right| dt \rightarrow 0 \text{ a.s.} \\
 & \leq \int_{\mathbb{R}} \left| \pi^*(t | x_1, \dots, x_n) - \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2} t^2 I(\theta_0)} \right| dt \\
 & \quad + \underbrace{\int_{\mathbb{R}} \left| \frac{\sqrt{I(\theta_0)}}{\sqrt{2\pi}} e^{-\frac{1}{2} t^2 I(\theta_0)} - \frac{\sqrt{\hat{I}_n}}{\sqrt{2\pi}} e^{-\frac{1}{2} t^2 \hat{I}_n} \right| dt}_{\rightarrow 0 \text{ by DCT}} \quad (\rightarrow 0, \text{ 1st part})
 \end{aligned}$$

Here we completed the proof of BYM Theorem.  
 The following theorem states that in the regular case, with the large sample, a Bayes estimate is approximately same as the MLE  $\hat{\theta}_n$ .

If we consider the squared error loss, the Bayes estimate for  $\theta$  is given by the posterior mean

$$\theta_n^* = \int \theta \pi_n(\theta | x_1, \dots, x_n) d\theta$$

- Theorem:- In addition to the assumption of the previous theorem, assume that  $\int |\theta| \pi(\theta) d\theta < \infty$ .  
 Then  $\sqrt{n}(\theta_n^* - \hat{\theta}_n) \rightarrow 0$  with probability one.

Proof:- Proceeding as in the proof of the previous theorem and using the assumption of finite mean of  $\pi(\cdot)$  we can prove

$$\int_{\mathbb{R}} |t| \left| \pi_n^*(t | x_1, \dots, x_n) - \sqrt{\frac{I(\theta_0)}{2\pi}} e^{-\frac{1}{2} t^2 I(\theta_0)} \right| dt \rightarrow 0 \text{ a.s.} \quad [\text{CHECK}]$$

This implies  $\int_{\mathbb{R}} t \pi_n^*(t | x_1, \dots, x_n) dt \rightarrow \int_{\mathbb{R}} t \sqrt{\frac{I(\theta_0)}{2\pi}} e^{-\frac{1}{2} t^2 I(\theta_0)} dt$

$$\begin{aligned}
 \text{Now, } \theta_n^* &= E(\theta | x_1, \dots, x_n) = 0 \text{ a.s.} \\
 &= E\left(\hat{\theta}_n + \frac{t}{\sqrt{n}} | x_1, \dots, x_n\right) = \hat{\theta}_n + \frac{1}{\sqrt{n}} \int t \pi_n^*(t | x_1, \dots, x_n) dt
 \end{aligned}$$

$$\therefore \sqrt{n}(\theta_n^* - \hat{\theta}_n) = \int t \pi_n^*(t | x_1, \dots, x_n) dt \rightarrow 0 \text{ a.s.}$$

Remark:-  $\sqrt{n}(\theta_n^* - \theta_0) \xrightarrow{d} N(0, I^{-1}(\theta_0))$  [Sampling Theory Result].

Suppose  $X \sim P$  and  $P \in \mathcal{P}$

Laplace has given the following idea of consistency:  
 If  $X_1, X_2, \dots$  iid Bernoulli( $\theta$ ) and  $\pi(\theta)$  is prior density that is continuous and (+ve) on  $(0,1)$ , then the posterior is consistent at all  $\theta_0 \in (0,1)$ .

An elementary proof of Laplace's result for a Beta Prior is as follows:

$$\pi(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\theta | X_1, \dots, X_n \sim \text{Beta}\left(\sum_{i=1}^n X_i + \alpha, n - \sum_{i=1}^n X_i + \beta\right)$$

$$E(\theta | X_1, \dots, X_n) = \frac{\sum_{i=1}^n X_i + \alpha}{\alpha + \beta + n} \rightarrow \theta_0 \text{ as } n \rightarrow \infty.$$

$$\text{var}(\theta | X_1, \dots, X_n) = \frac{\left(\sum_{i=1}^n X_i + \alpha\right) \left(n - \sum_{i=1}^n X_i + \beta\right)}{(\alpha + \beta + n)^2 (\alpha + \beta + n + 1)} \rightarrow 0 \text{ as } n \rightarrow \infty$$

In most unusual examples, posterior distr.s are consistent. We consider now, the interesting situations when

$$\Theta = \{\theta_1, \dots, \theta_k\}, \text{ a finite set.}$$

Theorem:- Let  $\Theta = \{\theta_1, \dots, \theta_k\}$  and  $X_1, X_2, \dots$  be iid observations with common density  $f(x|\theta), \theta \in \Theta$ . Consider a prior  $\{p_1, \dots, p_k\}$  where  $p_i > 0 \forall i=1(1)k$  and  $\sum_{i=1}^k p_i = 1$  and  $p_i = \text{Pr.}(\theta = \theta_i)$ .

Suppose  $\theta_t \in \Theta$  is the true value of  $\theta$  and  $\theta_i$ 's are distinguishable  $\theta \forall i \neq t$

$$\int f(x|\theta_t) \log\left(\frac{f(x|\theta_t)}{f(x|\theta_i)}\right) dx > 0.$$

Then w.p. 1

$$\lim_{n \rightarrow \infty} p(\theta_t | X_1, \dots, X_n) = 1 \text{ and}$$

$$\lim_{n \rightarrow \infty} p(\theta_i | X_1, \dots, X_n) = 0, \quad i \neq t.$$

Define,  $I(\theta) = E \left( \frac{\partial}{\partial \theta} \log f(x, \theta) \right)^2$

If  $I(\theta) = 0 \Rightarrow \frac{\partial}{\partial \theta} \log f(x, \theta)$  is small.

If  $I(\theta)$  is large,  $\Rightarrow \frac{\partial}{\partial \theta} \log f(x, \theta)$  is large.

Hence when  $I(\theta)$  is small, the rate of change of likelihood w.r.t.  $\theta$  remains constant.

Proof:-

$$p(\theta_i | x_1, \dots, x_n) = \frac{p_i f(x_1, \dots, x_n | \theta_i)}{\sum_{j=1}^k p_j f(x_1, \dots, x_n | \theta_j)}$$

$$= \frac{p_i \prod_{n=1}^n f(x_n | \theta_i)}{\sum_{j=1}^k p_j \prod_{n=1}^n f(x_n | \theta_j)}$$

To show  $p(\theta_i | x_1, \dots, x_n) \rightarrow 1$  if  $\theta_i = \theta_t$   
 $\rightarrow 0$  if  $\theta_i \neq \theta_t$

$$\therefore p(\theta_i | x_1, \dots, x_n) = \frac{p_i \prod_{n=1}^n \frac{f(x_n | \theta_i)}{f(x_n | \theta_t)}}{\sum_{j=1}^k p_j \prod_{n=1}^n \frac{f(x_n | \theta_j)}{f(x_n | \theta_t)}}$$

$$= \frac{p_i e^{n s_i}}{\sum_{j=1}^k p_j e^{n s_j}}, \text{ where } s_i = \frac{1}{n} \sum_{n=1}^n \log \frac{f(x_n | \theta_i)}{f(x_n | \theta_t)}$$

$$= \frac{p_i e^{n s_i}}{p_t + \sum_{1 \leq j \neq t \leq k} p_j e^{n s_j}}$$

The rest of the proof is skipped.

Doob's Theorem:- Suppose  $(H)$  and  $\mathcal{X}$  are both complete separable metric spaces (or, polish spaces) endowed with their respective Borel  $\sigma$ -fields  $\mathcal{B}(H)$  and  $\mathcal{A}$  and let  $\theta \rightarrow P_\theta$  be one to one.

Let,  $\pi$  be a prior and  $\{\pi(\cdot | x_1, \dots, x_n)\}$  be a posteriors, then  $\exists (H)_0 \subset (H)$  with  $\pi((H)_0) = 1$  s.t.  $\{\pi(\cdot | x_1, \dots, x_n)\}_{n \geq 1}$  is consistent at every  $\theta \in (H)_0$ .

$$P_{\theta_0}^{\infty} \left( \sqrt{n}(\theta - \theta_0) \right) \underset{\mathcal{L}}{\xrightarrow{d}} N \left( 0, \frac{1}{I(\theta_0)} \right)$$

Bayesian Information Criterion (BIC):-

Consider a model with likelihood  $f(x|\underline{\theta})$  and prior  $\pi(\underline{\theta})$ .

Then the integrated likelihood

$$m(\underline{x}) = \int q(\theta) e^{nh(\theta)} d\theta,$$

where,  $q(\theta) = \pi(\underline{\theta})$  and  $nh(\underline{\theta}) = \log f(\underline{x}|\underline{\theta})$ .

The Laplace approximation formula can be used to find an approximation to  $m(\underline{x})$ .

Schwarz (1978) proposed a criterion, known as the BIC based on this approximation, ignoring terms that stay bounded as the sample size  $n \rightarrow \infty$ .

The criterion is given by,

$$\log \hat{m}(\underline{x}) = \text{BIC} = \log f(\underline{x}|\hat{\theta}) - \frac{p}{2} \log n$$

Note that this approximation is free from the choice of prior.

$$\hat{m}(\underline{x}) = \frac{e^{\log f(\underline{x}|\hat{\theta})} \pi(\hat{\theta}) (\sqrt{2\pi})^p}{\sqrt{|-n^p C(\hat{\theta})|}}$$

$$\log \hat{m}(\underline{x}) = \log f(\underline{x}|\hat{\theta}) - \frac{p}{2} \log n.$$

# BAYESIAN HYPOTHESIS TESTING & MODEL SELECTION

On the irreconcilability of the p-value and Bayesian measures of evidence - an example.

$$X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$$

To test,  $H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$

We reject in frequentist inference if  $T(\underline{x}) = |\sqrt{n}(\bar{x} - \theta_0)|$  is p-value =  $P_{\theta_0} [|\sqrt{n}(\bar{x} - \theta_0)| > t]$ , where,  $t$  is the observed value of  $T(\underline{x}) = |\sqrt{n}(\bar{x} - \theta_0)|$

$$\text{i.e., } t = |\sqrt{n}(\bar{x} - \theta_0)|$$

Assume,  $P(H_0) = P(H_1) = 1/2$ .

Consider a prior degenerate at  $\theta_0$  and  $H_0$  and  $N(\mu, \tau^2)$  priors for  $\theta$  under  $H_1$ .

$$\pi(\theta) = \frac{1}{2} I(\theta = \theta_0) + \frac{1}{2} N(\mu, \tau^2) I(\theta \neq \theta_0)$$

$B_{01}$  = Bayes factor of  $H_0$  to  $H_1$

$$= \frac{f(\underline{x} | \theta_0)}{\int f(\underline{x} | \theta) \pi(\theta) d\theta}$$

Verify:

$$\sqrt{1 + \frac{1}{\rho^2}} \exp \left[ -\frac{1}{2} \left\{ \frac{(t - \rho\eta)^2}{1 + \rho^2} - \eta^2 \right\} \right]$$

Taking  $\mu = \theta_0$ ,  $\tau = 1$ , we have

$$B_{01} = \sqrt{1 + n} \exp \left\{ -\frac{1}{2} \cdot \frac{t^2}{1 + \frac{1}{n}} \right\}$$

$$\text{where, } \rho = \frac{1}{\sqrt{n}\tau}, \eta = \frac{(\theta_0 - \mu)}{\tau}$$

Take  $n = 50$ ,  $t = 1.96$

$$p\text{-value} = 0.05, B_{01} = 1.08$$

$$P(H_0 | \underline{x}) = 0.52$$

On the comparison of p-value and posterior probabilities:

$$X_1, \dots, X_n | \theta \stackrel{iid}{\sim} N(\theta, 1)$$

$H_0: \theta = \theta_0$  vs.  $H_1: \theta \neq \theta_0$

Notation:-  $B$  = Bayes factor of  $H_0$  relative to  $H_1$

$$\rho = P(H_0 | \underline{x})$$

$n$  = sample size

t	$\alpha$	1		5		10		20		50		100	
		B	P	B	P	B	P	B	P	B	P	B	P
1.645	0.10	0.72	0.42	0.79	0.44	0.89	0.47	1.27	0.56	1.86	0.64	...	...
1.960	0.05	0.54	0.35	0.49	0.33	0.59	0.37	1.08	0.52	1.08	0.52	-	-
2.576	0.01	0.27	0.21	0.15	0.13	0.16	0.14	0.28	0.22	0.28	0.22	-	-
3.291	0.001	0.10	0.09	0.03	0.03	0.02	0.02	0.03	0.03	0.03	0.03	-	-

Remarks:-

1. To a Bayesian the posterior probability of  $H_0$  summarises the evidence against  $H_0$ . In many common cases of testing  $H_0$ , p-value is smaller than the posterior prob. of  $H_0$ , by an order of magnitude. It may be noted that p-value ignores the likelihood of the data under the alternative and takes into account not only the observed deviation of the data from null hypothesis, as measured by test statistic but also more significant deviation.
2. One may try to see what happens when one changes the prior density of  $\theta$  under  $H_1$ . We will see that the difference between p-value and posterior probabilities persists even when we consider lower bounds on  $B$  and  $P$  over large class of priors under the alternative.

Note: We note that the differences between p-values and the corresponding Bayesian measure of evidence remain irreconcilable even when the lower bounds on such measures are considered. In other words, the least possible Bayes factor and posterior prob. of  $H_0$  are substantially larger than the corresponding p-values.



Suppose we are comparing  $k$  competing models:  $M_1, M_2, \dots, M_k$ .  
 for data  $X$  such that  $M_i: X$  has density  $f_i(x|\theta_i)$   
 for  $i=1(1)k$ . Here  $\theta_i$  are the unknown model parameters.  
 We will use hypothesis testing and model selection almost  
 interchangeably in this discussion.

Eg: The hypothesis testing problem. Suppose  $X_1, \dots, X_n$  have joint  
 densities  $f(x|\theta)$ , where  $\theta \in \mathcal{H}$ ,  $\mathcal{H}$  being the model space.

Example 1:  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(\theta, 1)$ ,  $\theta \in \mathbb{R}$

Suppose we want to test  $H_0: \theta \leq \theta_0$  Vs.  $H_1: \theta > \theta_0$ .

$H_0$  and  $H_1$  induce a partition of the parameter space as  
 $\mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1$ , where  $\mathcal{H}_0 = \{\theta \leq \theta_0\}$  and  $\mathcal{H}_1 = \{\theta > \theta_0\}$ .

This can be framed also as a model selection problem where  
 we want to choose between:

$$M_0: X_1, \dots, X_n | \theta \stackrel{iid}{\sim} f(x|\theta) \text{ with } \theta \in \mathcal{H}_0$$

$$M_1: X_1, \dots, X_n | \theta \stackrel{iid}{\sim} f(x|\theta) \text{ with } \theta \in \mathcal{H}_1$$

Example 2: Regression Models:

$Y$  is the "dependent" variable and  $x_1, x_2, \dots, x_p$  are the  
 "potential" regressors. If we use all the regressors the "full"  
 linear regression model is:

$$Y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} + \epsilon_i \text{ where } \epsilon_i \stackrel{iid}{\sim} N(0, 1)$$

One possible class of models is the nested class.

$$M_0: \beta_0 \in \mathbb{R}, \beta_i = 0 \quad \forall i=1(1)p$$

$$M_1: (\beta_0, \beta_1) \in \mathbb{R}^2, \beta_i = 0 \quad \forall i=2(1)p$$

$$M_2: (\beta_0, \beta_1, \beta_2) \in \mathbb{R}^3, \beta_i = 0 \quad \forall i=3(1)p$$

$$M_p: (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1} \text{ (the full model)}$$

$$M_0 \subset M_1 \subset \dots \subset M_p$$

Example 3:  $X \sim f(x|\theta), \theta \in \mathbb{R}$   
 $M_0: f(x|\theta) = \text{Cauchy}(\theta, 1)$   
 $M_1: f(x|\theta) = \text{Normal}(\theta, 1).$

Suppose given models:  $M_1, M_2, \dots, M_k$ . We have prior distr.  $\pi_i(\theta_i)$ ,  $i=1(1)k$  for the unknown parameters. Define the marginal distr. or predictive distr. of  $X$  given by

$$m_i(\underline{x}) = \int f_i(\underline{x}|\theta_i) \pi_i(\theta_i) d\theta_i, \quad i=1(1)k.$$

The Bayes factor of  $M_j$  to  $M_i$  is given by

$$B_{ji}(\underline{x}) = \frac{m_j(\underline{x})}{m_i(\underline{x})} = \frac{\int f_j(\underline{x}|\theta_j) \pi_j(\theta_j) d\theta_j}{\int f_i(\underline{x}|\theta_i) \pi_i(\theta_i) d\theta_i}$$

The Bayes factor is often interpreted as the odds provided by the data for  $M_j$  in comparison to  $M_i$ . Thus  $B_{ji} = 10$  would suggest that the data favour  $M_j$  over  $M_i$  at odds of 10 to one.

Alternatively,  $B_{ji}$  may also be thought of weighted likelihood ratio of  $M_j$  to  $M_i$ , the priors being the weight function. If prior probabilities  $P(M_j)$ ,  $j=1(1)k$  of the models are available one can compute the posterior probability of the models from the Bayes factors. It is easy to see that the posterior prob. of  $M_i$  given  $\underline{x}$  is

$$P(M_i|\underline{x}) = \frac{P(M_i) m_i(\underline{x})}{\sum_{j=1}^k P(M_j) m_j(\underline{x})} = \left[ \sum_{j=1}^k \frac{P(M_j)}{P(M_i)} B_{ji} \right]^{-1} \quad *$$

A particularly common choice (although not necessarily the best option in complex problem) is

$$P(M_i) = \frac{1}{k} \text{ for } i = 1(1)k.$$

In this case 
$$P(M_i | \underline{x}) = \frac{m_i(\underline{x})}{\sum_{j=1}^k m_j(\underline{x})} = \bar{m}_i(\underline{x}).$$

In scientific reporting it is common to provide  $\bar{m}_i(\underline{x})$  rather than  $P(M_j | \underline{x})$  since the prior probabilities of models can be a contentious matter and anyone can use the  $\bar{m}_i(\underline{x})$  to determine their personal posterior probabilities via (\*), noting that

$$B_{ji} = \frac{\bar{m}_j(\underline{x})}{\bar{m}_i(\underline{x})}$$

Simple Case:  $k = 2$

$$B_{10} = \frac{m_1(\underline{x})}{m_0(\underline{x})}$$

$$P(M_0) = \pi_0 \quad P(M_1) = 1 - \pi_0$$

$$P(M_1 | \underline{x}) = \frac{P(M_1) m_1(\underline{x})}{\sum_{j=0}^k P(M_j) m_j(\underline{x})}$$

$$\frac{P(M_1 | \underline{x})}{P(M_0 | \underline{x})} = B_{10} \frac{P(M_1)}{P(M_0)}$$

$$P(M_0 | \underline{x}) = \frac{P(M_0) m_0(\underline{x})}{\sum_{j=0}^k P(M_j) m_j(\underline{x})}$$

i.e.,  $B_{10} = \frac{\text{posterior odds}}{\text{prior odds}}$

$$\frac{c_1 \int f_i(\underline{x} | \theta_i) \pi_i^N(\theta_i) d\theta_i}{c_2 \int f_j(\underline{x} | \theta_j) \pi_j^N(\theta_j) d\theta_j}$$

measure of evidence changes for same data, So, using non-informative prior doesn't justify the problem.

## Motivation for Bayesian approach to model selection:

1. Bayes factors and posterior model probability are easy to understand. The interpretation of Bayes factors are easily understandable by non-statisticians.
2. Bayesian model selection is consistent under mild condition in finite/fixed dimensional problem.
3. Bayesian procedures naturally penalize model complexity and they do need an explicit introduction of the penalty term. In this way, simpler models are preferred to complex models.
4. The Bayesian approach to model selection is conceptually the same, regardless of the no. of models under consideration.

Example:  $H_0: \theta \in \mathbb{H}_0$  Vs.  $H_1: \theta \in \mathbb{H} - \mathbb{H}_0 = \mathbb{H}_1$ .

Case 1:  $\mathbb{H}_0$  and  $\mathbb{H}_1$  are of same "dimension"

(viz.  $H_0: \theta \leq \theta_0$  Vs.  $H_1: \theta > \theta_0$  where  $X_1, \dots, X_n | \theta \stackrel{iid}{\sim} N(\theta, 1)$ ,  $\theta \sim N(\mu, \tau^2)$ )  
 We will calculate  $P(H_0 | \tilde{x})$  and  $P(H_1 | \tilde{x})$  and check which one is larger.

$$\pi(\theta \leq \theta_0) = \int \pi(\theta) d\theta = \pi_0 \quad \text{and} \quad \pi(\theta > \theta_0) = 1 - \pi_0.$$

$$\begin{aligned} \pi(\theta) &= P(H_0) P(\theta | H_0) + P(H_1) P(\theta | H_1) \\ &= \pi_0 \cdot \frac{N(\mu, \tau^2)}{\pi_0} + (1 - \pi_0) \frac{N(\mu, \tau^2)}{(1 - \pi_0)} \end{aligned}$$

$$= \pi_0 \pi_1(\theta) \mathbb{I}(\theta \in \mathbb{H}_0) + (1 - \pi_0) \pi_2(\theta) \mathbb{I}(\theta \in \mathbb{H}_1), \text{ where}$$

$$\pi_1(\theta) = \frac{\exp\left(-\frac{1}{2\tau^2}(\theta - \mu)^2\right)}{\sqrt{2\pi} \tau \pi_0}$$

Case 2:  $\mathbb{H}_0$  and  $\mathbb{H}_1$  are not of same dimension, viz.  
 $H_0: \theta = \theta_0$  Vs.  $H_1: \theta \neq \theta_0$ , where,  
 $X_1, \dots, X_n | \theta \stackrel{iid}{\sim} N(\theta, 1)$

$$\pi(\theta) = \pi_0 \mathbb{I}(\theta = \theta_0) + (1 - \pi_0) \pi_1(\theta) \mathbb{I}(\theta \neq \theta_0)$$

## Motivation behind objective Bayesian model selection; —

There has been a long debate in the Bayesian community about the roles of subjective and objective Bayesian analysis. Subjective Bayesian analysis is an attractive proposition but few would disagree that it is frequently not a realistic possibility because it is sometimes simply not possible to obtain the extensive needed elicitation from subjective experts.

In model selection, this last argument is particularly compelling because we often initially entertain a wide variety of models and careful subjective specification of prior distributions for all the parameters of the model is essentially impossible. Four methods for developing default Bayesian model selection have come in modern research. There are the conventional priors approach, the BIC approach, the Intrinsic Bayes factor (IBF) approach, the fractional Bayes Factor approach.

Conventional priors approach: -  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ ;  $\mu, \sigma^2$  unknown.

to test  $H_0: \mu = 0$  vs.  $H_1: \mu \neq 0$

Jeffrey's suggestion: Take prior  $g_0(\sigma) = \frac{1}{\sigma}$  for  $\sigma$  under  $H_0$ . Under  $H_1$ , take the same prior for  $\sigma$  and take the conventional priors for  $\mu$  and  $\sigma$  as  $g_1(\mu|\sigma) = \frac{1}{\sigma} g_2(\mu|\sigma)$ , where  $g_2(\cdot)$  is a pdf. Then under  $H_1$  the prior is

$$g_1(\mu, \sigma) = \frac{1}{\sigma} g_1(\mu|\sigma)$$

An initial choice of  $g_2(\cdot)$  is  $N(0, \tau^2)$ , one can take  $c=1$ . Jeffrey noted that for reasonable prior specification  $BF_{01} \rightarrow 0$  if  $\bar{x} \rightarrow \infty$  and  $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$  is bounded.

He gave an argument which implies that unless  $g_1$  has no finite moments this will not happen. In particular,

$g_2 = \text{Normal}$   $BF_{01} \not\rightarrow 0$  under above criterion.

Jeffrey suggested to use

$$g_2 = \text{cauchy}$$

So, Jeffrey's final recommendation are

$$g_0(\sigma) = \frac{1}{\sigma} \text{ under } H_0$$

$$g_1(\mu, \sigma) = \frac{1}{\sigma} g_2(\mu/\sigma) = \frac{1}{\sigma} \cdot \frac{1}{\sigma \pi \left(1 + \frac{\mu^2}{\sigma^2}\right)} \text{ under } H_1.$$

Jeffrey's test in normal problems:

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$$H_0: \mu = 0 \quad \text{Vs.} \quad H_1: \mu \neq 0$$

$$\text{Prior suggested: } g_0(\sigma) = \frac{1}{\sigma} \text{ under } H_0$$

$$g_1(\mu, \sigma) = \frac{1}{\sigma} \cdot \frac{1}{\sigma \pi \left(1 + \frac{\mu^2}{\sigma^2}\right)} \text{ under } H_1.$$

$$BF_{01} = \frac{\int f(x_1, \dots, x_n | 0, \sigma^2) g_0(\sigma) d\sigma}{\iint f(x_1, \dots, x_n | \mu, \sigma^2) g_1(\mu, \sigma) d\mu d\sigma}$$

There is no closed form for the expression for the denominator of  $BF_{01}$ . To deal with it one may note that

$$g_1(\mu | \sigma^2) = \int_0^{\infty} \frac{\sqrt{\lambda}}{\sqrt{2\pi} \sigma} \exp\left(-\frac{\lambda}{2\sigma^2} \mu^2\right) \frac{1}{\sqrt{2\pi}} \lambda^{\frac{1}{2}-1} e^{-\lambda/2} d\lambda.$$

$$\text{i.e., } \mu | \lambda \sim N(0, \sigma^2/\lambda), \quad \lambda \sim \text{Gamma}(\frac{1}{2}, \frac{1}{2}).$$

To calculate the denominator of  $BF_{01}$ , one can integrate over  $\mu$  and  $\sigma$  in closed form using the mixture representation. Ultimately, one has a one-dimensional integral over  $\lambda$  left which is easier to handle.

A Curious Example:- Einstein's theory of gravitation predicts deflection of light due to gravitation and specifies the amount of deflection. He predicted that light of stars would deflect stars, under gravitational pull of the sun on the nearby stars, but the effect would be invisible only during a total solar eclipse when the deflection can be measured through the apparent change in stars' position.

A famous experiment by a team lead by Eddington, led to acceptance of Einstein's theory. There are four observations, two collected in 1919 in Eddington's expedition and two more by other groups by in 1922 and 1929. These are  $\alpha_1 = 1.98$ ,  $\alpha_2 = 1.61$ ,  $\alpha_3 = 1.18$  and  $\alpha_4 = 2.24$ .

Suppose they are iid  $N(\mu, \sigma^2)$ , both parameters unknown. Einstein's theory says  $\mu$  should be 1.75.

We test  $H_0: \mu = 1.75$  Vs.  $H_1: \mu \neq 1.75$

If we use the conventional priors of Jeffreys then BF01 turns out to be 2.98.

A common choice of prior for the normal linear model is the conjugate prior, called a  $g$ -prior in Zellner (1986).

For a linear model  $\mu: y = X\beta + \epsilon$ ,  $\epsilon \sim N_n(0, \sigma^2 I_n)$

where,  $\sigma^2$  and  $\beta = (\beta_1, \dots, \beta_k)$  are unknown and  $X$  is  $(n \times k)$  given design matrix of rank  $k < n$ , the  $g$ -prior density is given by

$$\pi(\sigma) = \frac{1}{\sigma}, \quad \pi(\beta | \sigma) \sim N_k \left( 0, g\sigma^2 (X'X)^{-1} \right)$$

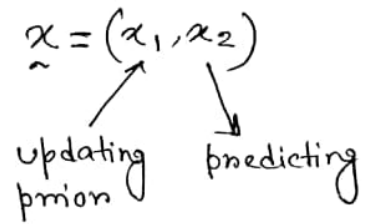
The key advantage of  $g$ -prior is that the marginal density,  $m(\underline{y})$ , is available in closed form and is given by,

$$m(\underline{y}) = \frac{|\underline{y}'\underline{y}|^{n/2}}{2\pi^{n/2} (1+g)^{k/2}} \left( \underline{y}'\underline{y} - \frac{g}{1+g} \underline{y}'X(X'X)^{-1}X'\underline{y} \right)^{-n/2}$$

Thus BF and posterior model probability for comparing any two linear models are available.

$$BF_{01}(\tilde{x}) = \frac{\int f_1(\tilde{x} | \theta_1) \pi_1(\theta_1) d\theta_1}{\int f_2(\tilde{x} | \theta_2) \pi_2(\theta_2) d\theta_2}$$

$$BF_{01}(\tilde{x}_2) = \frac{\int f_1(x_2 | \theta_1) \pi(\theta_1 | x_1) d\theta_1}{\int f_2(x_2 | \theta_2) \pi(\theta_2 | x_1) d\theta_2}$$



The Intrinsic Bayes Factor Approach: Consider two models  $M_0$  and  $M_1$  for data  $X$

with density  $f_i(x | \theta_i)$  under  $M_i$ ,  $\theta_i$  being an unknown parameter of dimension  $p_i$ ,  $i=0,1$ . Given prior specification  $g_i(\theta_i)$  for parameter  $\theta_i$ , the Bayes factor of  $M_1$  to  $M_0$  is

$$B_{10} = \frac{m_1(\tilde{x})}{m_0(\tilde{x})} = \frac{\int f_1(x | \theta) g_1(\theta) d\theta}{\int f_0(x | \theta_0) g_0(\theta_0) d\theta_0} \dots \dots \dots (*)$$

When subjective specification of prior distr. is impossible, one would look for automatic method that uses standard non-informative priors that are typically improper.

If  $g_i$ 's are improper, they are defined only upto arbitrary multiplicative constant  $c_i$ . Then  $c_i g_i$  has as much validity as  $g_i$ .

This implies that  $(\frac{c_1}{c_0}) B_{10}$  has as much validity as  $B_{10}$ . Thus the Bayes factor is determined only upto arbitrary multiplicative constants.

Remark: It is a curious fact that the use of diffuse (flat) proper prior may not provide a good solution to the problem described about the use of proper prior. Truncation of non-informative priors leads to large penalty for the more complex model.



## BAYESIAN COMPUTATION

Bayesian analysis requires computation of expectations and quantiles of probability distributions arising as posterior distributions. Sometimes modes of the posterior are also needed. If conjugate priors are not used, as in many cases, posterior distributions will not be standard distns. and hence the required Bayesian quantiles of interest can't be computed in closed form. So thus special techniques are needed for Bayesian computation.

Example: Suppose  $X \sim N(\theta, 1)$  and  $\theta \sim \text{Cauchy}(\mu, \tau)$  prior with  $\mu, \tau$  known. (say,  $\mu=0, \tau=1$ )

$$\pi(\theta|x) \propto \exp\left\{-\frac{(\theta-x)^2}{2}\right\} \left[\tau^2 + (\theta-\mu)^2\right]^{-1}$$

$$E(\theta|x) = \frac{\int_{-\infty}^{\infty} \theta \exp\left\{-\frac{(\theta-x)^2}{2}\right\} \left[\tau^2 + (\theta-\mu)^2\right]^{-1} d\theta}{\int_{-\infty}^{\infty} \exp\left\{-\frac{(\theta-x)^2}{2}\right\} \left[\tau^2 + (\theta-\mu)^2\right]^{-1} d\theta}$$

Easy to see that  $E(\theta|x)$  does not come out in any analytical tractable form. It is clear that some form of approximation will be needed. One may try numerical integration technique as the problem is low-dimensional.

Example: Suppose  $x_1, \dots, x_k$  are independent Poisson counts with  $x_i \sim \text{Poisson}(\theta_i)$ ,  $\theta_i$ 's are a-priori considered related, and a joint multivariate normal prior distns as their logarithm is assumed. Specifically, let  $\gamma_i = \log(\theta_i)$  and suppose  $\gamma = (\gamma_1, \dots, \gamma_k)$  has a distn  $\gamma \sim N_k\left(\mu, \frac{1}{\tau^2} \left\{ (1-\rho)I_k + \rho \frac{11'}{k} \right\}\right)$  where,  $\frac{1}{\tau^2}$  is the k-vector with all element being 1 ( $\mu, \tau^2, \rho$  are known).

$$\pi(\gamma|x) \propto \exp\left\{-\sum_{i=1}^k \{e^{\gamma_i} - \gamma_i x_i\} - \frac{1}{2\tau^2} (\gamma - \mu \frac{1}{k})' \left\{ (1-\rho)I_k + \rho \frac{11'}{k} \right\}^{-1} (\gamma - \mu \frac{1}{k})\right\}$$

$$P(x|\theta) = \frac{e^{-\theta} \theta^x}{x!} = e^{-\theta} \frac{\theta^x}{e^{\theta \log \theta}}$$

$$E(\theta_j | \tilde{x}) = E(\exp(\gamma_j) | \tilde{x})$$

$$= \frac{\int_{\mathbb{R}^k} \exp(\gamma_j) g(\tilde{\gamma} | \tilde{x}) d\tilde{\gamma}}{\int_{\mathbb{R}^k} g(\tilde{\gamma} | \tilde{x}) d\tilde{\gamma}}$$

(Curse of dimensionality) The above is a ratio of two  $k$ -dimensional integrals and as  $k$  grows, the integrals become less easy to work with. Numerical integration fails to be an efficient option in that case.

The errors in approximation associated with the numerical methods increase as a power of the dimensions in numerical integration may not be a good option except in very very low dimensional problem.

The way out of these problems, in many cases is through other approximation techniques. One could try, in some cases, the Laplace approximation, EM algorithm etc. But here we will stress on general purpose simulation based techniques.

Monte Carlo Sampling:- An alternative to numerical integration or analytical approximation to compute an integral is statistical sampling. This probabilistic technique is a familiar tool in statistical inference. To estimate a population mean or a population proportion, a natural approach is to gather a large sample from this population and to consider the corresponding sample mean or sample proportion. The law of large numbers guarantee that the estimates thus obtained are good provided the sample size is large.

Ref: "Markov chain Monte Carlo in practice" by Spigethalter et al.

Suppose  $f$  be a density (or pmf) and we want to estimate

$$E_f h(x) = \int h(x) f(x) dx.$$

If iid observation  $x_1, x_2, \dots$  can be generated from  $f$ , then

$$\bar{h}_m = \frac{1}{m} \sum_{i=1}^m h(x_i) \xrightarrow[\text{a.s.}]{P} E_f h(x).$$

This provides a justification of using  $\bar{h}_m$  as an approximation for  $E_f h(x)$  for large  $m$ .

Example (Continued)

$$E(\theta|x) = \frac{\int_{-\infty}^{\infty} \theta \phi(\theta-x) (\tau^2 + (\theta-\mu)^2)^{-1} d\theta}{\int_{-\infty}^{\infty} \phi(\theta-x) (\tau^2 + (\theta-\mu)^2)^{-1} d\theta};$$

where  $\phi(\cdot)$  denotes the density of  $N(0,1)$ .

So,  $E(\theta|x)$  is the ratio of expectations,  $h_1(\theta) = \frac{\theta}{\tau^2 + (\theta-\mu)^2}$  to that of  $h_2(\theta) = \frac{1}{\tau^2 + (\theta-\mu)^2}$  w.r.t.  $N(x,1)$  distn.

So, we may simply sample  $\theta_1, \theta_2, \dots$  from  $N(x,1)$  and we have

$$\hat{E}(\theta|x) = \frac{\sum_{i=1}^m \theta_i (\tau^2 + (\theta_i - \mu)^2)^{-1}}{\sum_{i=1}^m (\tau^2 + (\theta_i - \mu)^2)^{-1}}$$

as our Monte Carlo estimate of  $E(\theta|x)$ .

Another approach could be to write  $E(\theta|x)$  as

$$E(\theta|x) = \frac{\int_{-\infty}^{\infty} \theta \exp(-(\theta-x)^2) \pi(\theta) d\theta}{\int_{-\infty}^{\infty} \exp(-(\theta-x)^2) \pi(\theta) d\theta}$$

when  $\pi(\theta) = \text{Cauchy}(\mu, \tau)$ .

A second Monte Carlo estimate is given by

$$\hat{E}(\theta|x) = \frac{\sum_{i=1}^m \theta_i \exp(-(\theta_i-x)^2)}{\sum_{i=1}^m \exp(-(\theta_i-x)^2)}$$

But one has to note that in either approach one doesn't necessarily get a sample, that would closely resemble a bonafide sample from the actual posterior distn.

Thus it is of need to find techniques of approximately generating samples from a distn, which is either the posterior distn itself or a distn, resembling the posterior density closely.

## Markov Chain Monte Carlo (MCMC):-

MCMC refers to a broad class of iterative procedure that produce random sequence with the Markov properties such that the Markov chain has the posterior distr., as its stationary distr. and the chain converges to that distr. So, the technique is to produce a sequence  $X_t$  of RVs with such property. Lets first do a quick review of MC.

A prob. distr.  $\pi$  is called stationary or invariant for a transition prob.  $P$  on the associated MC  $\{X_n\}$  if it is the case that when the prob. distr. of  $X_0$  is  $\pi$  then the same is true for  $X_n$  for all  $n \geq 1$ .

Thus, in the countable <sup>state</sup> space case a prob. distr.  $\pi = \{\pi_i : i \in S\}$  is stationary for  $P$ , if  $\forall j \in S$ ,

$$\begin{aligned} P(X_1 = j) &= \sum_i P(X_1 = j | X_0 = i) P(X_0 = i) \\ &= \sum_i \pi_i p_{ij} = P(X_0 = j) = \pi_j. \end{aligned}$$

In vector notation,  $\tilde{\pi} = (\pi_1, \pi_2, \dots)$  satisfies  $\tilde{\pi} = \tilde{\pi} P$ .

A Markov chain is called irreducible if for any two states  $i$  and  $j$ ,  $P(X_n = j | X_0 = i) > 0$  for some  $n \geq 1$ .

Theorem:- Let  $\{X_n\}_{n \geq 0}$  be a MC with a countable state space  $S$  and a transition probability matrix  $P$ . Further suppose, it is irreducible with stationary prob. distr.  $\pi = \{\pi_i : i \in S\}$ . Then for any odd functn.  $h: S \rightarrow \mathbb{R}$  and for any initial distr. of  $X_0$ ,

$$\frac{1}{n} \sum_{i=0}^{n-1} h(X_i) \rightarrow \sum_j h(j) \pi_j \text{ in probability as } n \rightarrow \infty.$$

To see how this is useful to us, note that given a prob. distr.  $\pi$  on  $S$  and a function  $h$  on  $S$ , suppose it is desired to find the "integral of  $h$  w.r.t.  $\pi$ " which reduces

to  $\sum_j h(j) \pi_j$  in the countable case.

Then look for an irreducible MC  $\{X_n\}$  with space  $S$  and stationary distr.  $\pi$ . Then starting from some initial value  $X_0$ , run the chain  $\{X_j\}$  for a period of time, say,  $0, 1, \dots, n-1$  and consider as an estimate

$$\mu_n = \frac{1}{n} \sum_{i=0}^{n-1} h(X_i) \text{ and by the theorem}$$

$$\mu_n \approx \sum_j h(j) \pi_j \text{ for large } n.$$

This technique is called MCMC method.

An irreducible MC  $\{X_n\}$  with a countable state space  $S$  is called aperiodic if for some  $i \in S$ , the gcd  $\{n : p_{ii}^{(n)} > 0\} = 1$ . Then, in addition to the theorem above, we have  $\sum_j |P(X_n = j) - \pi_j| \rightarrow 0$  as  $n \rightarrow \infty$ .

for any initial distr.  $X_0$ . In other words, for large  $n$ , the distr. of  $X_n$  will be close to  $\pi$ .

### Metropolis Hastings Algorithm:-

This is a very general MCMC method with wide applicability. The idea is not to directly simulate from given target density (which may be computationally very difficult) at all, but to simulate an easy Markov chain that has this target density as the density of its stationary distribution.

Let  $S$  be a finite (or countable) set. Let  $\pi$  be a prob. distr. of  $S$ . We shall call  $\pi$ , the target distribution. Let  $Q = ((q_{ij}))$  be a transition prob. matrix such that for each  $i$ , it is computationally easy to generate a sample from the distr.  $\{q_{ij} : j \in S\}$ .

Let us generate MC as follows:

If  $X_n = i$  first sample from the distribution.

$\{q_{ij} : j \in S\}$  and denote that observation as  $Y_n$ . Thus choose  $X_{n+1}$  from the two values  $X_n$  and  $Y_n$  accordingly to

$$P(X_{n+1} = Y_n | X_n, Y_n) = f(X_n, Y_n)$$

$$P(X_{n+1} = X_n | X_n, Y_n) = 1 - f(X_n, Y_n).$$

where the "acceptance probability"  $f(\cdot, \cdot)$  is given by

$$f(i, j) = \min \left\{ \frac{\pi_j}{\pi_i} \cdot \frac{q_{ji}}{q_{ij}}, 1 \right\} \text{ for all } (i, j) \text{ s.t. } \pi_i q_{ij} > 0.$$

Continuous Case: - Suppose  $S$  is a continuum and the target distribution  $\pi(\cdot)$  has the density  $p(\cdot)$ . Then MH proceeds as follows:

Let  $Q$  be a transition function s.t. for each  $x$ ,  $Q(x, \cdot)$  has density  $q(x, y)$ . Then proceed as in the discrete case. But set the "acceptance probability"  $f(x, y)$  to be

$$f(x, y) = \min \left\{ \frac{p(y)}{p(x)} \frac{q(y, x)}{q(x, y)}, 1 \right\}$$

for all  $(x, y)$  s.t.  $p(x)q(x, y) > 0$ .

Some popular choices of the proposal distn  $q(\cdot, \cdot)$  are as follows:  
In Metropolis et al. (1953),  $q(\cdot, \cdot)$  was such that  $q(x, y) = q(y, x) \forall (x, y)$ .

For example, when  $x$  is continuous  $q(x, \cdot)$  might be a multivariate normal with mean  $x$  and constant covariance matrix  $\Sigma$ .

For the symmetric proposal distn the acceptance probability

$$f(x, y) = \min \left\{ 1, \frac{\pi(y)}{\pi(x)} \right\}.$$

A special case of M-H algorithm is the random walk Metropolis for which  $q(x, y) = q(|x - y|)$ , while choosing a proposal distn, its scale need to be chosen carefully.

A continuous proposal distr. generating small steps  $Y - X_t$  will generally have a high acceptance rate but will nevertheless "mix" slowly i.e., more slowly around the support of the target  $\pi(\cdot)$ .

A bold proposal distribution generating large steps with the purpose moves from the body to the tails of the distr. giving small values of  $\frac{\pi(y)}{\pi(x_t)}$  and a low probability of acceptance. Such a chain will frequently not move, again resulting in slow mixing.

The independence sampler (Tierney, 1994) is a M-H algorithm where proposed  $q(x, y) = q(y)$  doesn't depend on  $x$ . For this the acceptance prob. can be written in the form  $\alpha(x, y) = \min\left(1, \frac{w(y)}{w(x)}\right)$ , where  $w(x) = \frac{\pi(x)}{q(x)}$ .

For the independence sampler to work well,  $q(\cdot)$  should be a good approximation of  $\pi(\cdot)$  but it is safest if  $q(\cdot)$  is heavier tailed than  $\pi(\cdot)$ . To see this, suppose  $q(\cdot)$  is heavier tailed than  $\pi(\cdot)$  and  $X_t$  is currently in the tails of  $\pi(\cdot)$ . Most candidates will not be in the tails so  $w(X_t)$  will be much larger than  $w(y)$  giving a low acceptance probability. The heavy tailed independent proposals help to avoid long periods stuck in the tails.

Gibbs Sampling:- The Gibbs sampler is a technique suitable for generating a Markov chain of nice properties that has its stationary distribution a target distribution in higher dimensional space. The most interesting aspect of this technique is that to run this MC, it suffices to generate observations from univariate distr.s. Suppose you want to simulate a MC with  $\pi(\theta | y)$  where  $\theta = (\theta_1, \dots, \theta_k)$  as its stationary distribution. In Gibbs sampling in order to do this you proceed as follows:

Start with  $\theta_0 = (\theta_{01}, \theta_{02}, \dots, \theta_{0k})$

Generate  $\theta_{11}$  from  $\pi_1(\theta_1 | \theta_{02}, \dots, \theta_{0k}, y)$

$\theta_{12}$  from  $\pi_2(\theta_2 | \theta_{01}, \theta_{03}, \dots, \theta_{0k}, y)$

$\vdots$   
 $\theta_{1k}$  from  $\pi_k(\theta_k | \theta_{01}, \theta_{02}, \dots, \theta_{0k-1}, y)$

This gives  $\theta_1 = (\theta_{11}, \theta_{12}, \dots, \theta_{1k})$

Then produce  $\theta_2$  from  $\pi_1(\theta_1 | \theta_{12}, \dots, \theta_{1k}, y)$  and so on.

Example:- (i) suppose  $\begin{pmatrix} x \\ y \end{pmatrix} \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$ .

$$x | y = y \sim N(\rho y, 1 - \rho^2)$$

$$y | x = x \sim N(\rho x, 1 - \rho^2)$$

To apply GS technique

1. Given  $x_i$  for  $x$ , draw  $y_i \sim N(\rho x_i, 1 - \rho^2)$
2. Given  $y_i$  for  $y$ , draw  $x_{i+1} \sim N(\rho y_i, 1 - \rho^2)$ .

(ii)  $x | \theta \sim N(\theta, \sigma^2)$ ;  $\sigma^2$  known.

$\theta \sim \text{cauchy}(\mu, \tau)$

$$\pi(\theta) \propto (\tau^2 + (\theta - \mu)^2)^{-1}$$

$$\propto \int_0^\infty \left(\frac{\lambda}{2\pi\tau^2}\right)^{1/2} \exp\left(-\frac{\lambda}{2\tau^2}(\theta - \mu)^2\right) \lambda^{1/2-1} \exp\left(-\frac{\lambda}{2}\right) d\lambda.$$



Thus  $\pi(\theta)$  may be considered the marginal prior density from the joint prior density of  $(\theta, \lambda)$ , where,  $\theta | \lambda \sim N(\mu, \sigma^2/\lambda)$ , and  $\lambda \sim \text{Gamma}(1/2)$ . It is clear that, this representation of  $\pi(\cdot)$  leads to an implicit hierarchical prior structure with  $\lambda$  as a hyper parameter.

Consequently,  $\pi(\theta | x)$  may be treated as the marginal density for  $\pi(\theta, \lambda | x)$ . To apply the GS note that

$$\theta | \lambda, x \sim N\left(\frac{\sigma^2}{\sigma^2 + \lambda \sigma^2} \bar{x} + \frac{\lambda \sigma^2}{\sigma^2 + \lambda \sigma^2} \mu, \frac{\sigma^2 \sigma^2}{\sigma^2 + \lambda \sigma^2}\right)$$

$$\lambda | \theta, x \sim \lambda | \theta \sim \exp\left(-\frac{\sigma^2 + (\theta - \mu)^2}{2\sigma^2}\right).$$

Example: Suppose we are studying the distribution of the number of defectives  $X$  in the daily productions of a product. Consider the model  $(x | y, \theta) \sim \text{Bin}(y, \theta)$ , where,  $y$ , a day's production, is a RV with Poisson distribution with unknown mean  $\lambda$  and  $\theta$  is the probability (unknown) that a product is defective.  $y$  is not observable and inference has to be made on the basis of  $X$  only. The prior distr. is such that

$$\theta | y = y \sim \text{Beta}(\alpha, \gamma) \text{ with known } \alpha \text{ and } \gamma \text{ independent of } y.$$

First note that  $x | \theta \sim \text{Poisson}(\lambda \theta)$  (check) since  $\theta \sim \text{Beta}(\alpha, \gamma)$ .

$$\text{We have } \pi(\theta | X = x) \propto \exp(-\lambda \theta) \theta^{x + \alpha - 1} (1 - \theta)^{\gamma - 1}; 0 < \theta < 1.$$

Note:  $\theta | x$  is not a standard distr. and hence posterior quantile can't be obtained in closed form. You can however use MCMC technique. Let's see how Gibbs sampling works here. Instead of focusing  $\theta | x$  directly view it as a marginal component of  $(y, \theta | x)$ . It is easy to see that the full conditionals of this are given by

$$y | x = x, \theta \sim x + \text{Poisson}(x(1 - \theta))$$

$$\text{and } \theta | x = x, y = y \sim \text{Beta}(\alpha + x, \gamma + y - x).$$

# HIERARCHICAL BAYES (HB) APPROACH AND PARAMETRIC EMPIRICAL BAYES (PEB) APPROACH

We consider  $p$  similar but not identical populations with densities  $f(x|\mu_1), \dots, f(x|\mu_p)$ . We illustrate below with normal densities:

Consider  $p$  independent random samples, each of size  $n$ , from  $p$  normal populations:

$$X_{ij}, i=1,2,\dots,n \text{ are iid } N(\mu_j, \sigma^2), j=1,2,\dots,p.$$

For simplicity, assume  $\sigma^2$  to be known. These  $p$  populations may correspond to  $p$  adjacent small areas with unknown per capita income  $\mu_1, \dots, \mu_p$ , as in small area estimation.

They could correspond to  $p$  clinical trials in a particular hospital with  $\mu_1, \dots, \mu_p$  being the mean effects of the drug being tested. It may be said that in these examples the different studied population are related to each other.

In order to assign a prior distribution for the  $\mu_j$ 's, we model them as exchangeable rather than i.i.d or just independent. An exchangeable, dependent structure is consistent with the assumption that the studies are similar in a broad sense, so they share many common elements. We use a standard way of generating exchangeable  $\mu_j$ 's which is as follows:

We introduce a hyperparameter  $\eta = (\eta_1, \eta_2)$  and assume  $\mu_1, \mu_2, \dots, \mu_p$  are iid  $N(\eta_1, \eta_2)$  given  $\eta$  and  $\eta = (\eta_1, \eta_2) \sim \pi(\cdot, \cdot)$  (some prior density for  $\eta$ )

$$\tilde{X} = (X_{ij}, i=1, \dots, n, j=1, \dots, p)$$
$$\tilde{\mu} = (\mu_1, \dots, \mu_p), \bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}$$

sufficient statistic  $\bar{X} = (\bar{X}_1, \dots, \bar{X}_p)$

Given  $\eta$

$$(\underline{x}, \underline{\mu}) \sim f(\underline{x}, \underline{\mu} | \eta) = \prod_{j=1}^p \left\{ \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_{ij} - \mu_j)^2\right] \frac{1}{\sqrt{2\pi}\eta_2} \exp\left[-\frac{1}{2\eta_2} (\mu_j - \eta_1)^2\right] \right\}$$

$$\pi(\underline{\mu} | \underline{x}, \eta) \propto \prod_{j=1}^p \left\{ \exp\left[-\frac{n}{2\sigma^2} (\mu_j - \bar{x}_j)^2\right] \cdot \exp\left[-\frac{1}{2\eta_2} (\mu_j - \eta_1)^2\right] \right\}$$

[Note that  $\sum_{i=1}^n (x_{ij} - \mu_j)^2 = \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 + n(\bar{x}_j - \mu_j)^2$ ]

Thus,  $(\mu_j | \underline{x}, \eta) \equiv (\mu_j | \bar{x}_j, \eta)$  are independent normal with

$$E(\mu_j | \bar{x}_j, \eta) = \frac{\frac{n}{\sigma^2} \bar{x}_j + \eta_1/\eta_2}{\frac{n}{\sigma^2} + \frac{1}{\eta_2}} = \frac{\eta_2 \bar{x}_j + (\sigma^2/n)\eta_1}{\eta_2 + (\sigma^2/n)}$$

$$= (1-B) \bar{x}_j + B \eta_1, \text{ where } B = \frac{\sigma^2/n}{\eta_2 + \frac{\sigma^2}{n}}$$

$$\text{Var}(\mu_j | \bar{x}_j, \eta) = \frac{1}{\frac{n}{\sigma^2} + \frac{1}{\eta_2}} = \frac{\eta_2 \sigma^2/n}{\eta_2 + \frac{\sigma^2}{n}}$$

Given  $\eta$ , Bayes estimate of  $\mu_j$  (w.r.t. square error loss) is

$$E(\mu_j | \underline{x}, \eta) = E(\mu_j | \bar{x}_j, \eta) = (1-B) \bar{x}_j + B \eta_1, \text{ where } B = \frac{\sigma^2/n}{\eta_2 + \sigma^2/n}$$

but  $\eta$  is unknown.

Hierarchical Bayes Approach (HB):-

Consider a prior for  $\eta = (\eta_1, \eta_2) \sim \pi(\cdot, \cdot)$  (prior density of  $\eta$ )  
Find the posterior density  $\pi(\eta | \underline{x})$  of  $\eta$ .

$(\underline{x}, \underline{\mu}) | \eta \sim f(\underline{x}, \underline{\mu} | \eta)$  as given above. Integrating out the  $\mu_j$ 's and holding  $\eta$  fixed, we get  $\bar{x}_j$ 's are independent and  $\bar{x}_j | \eta \sim N\left(\eta_1, \eta_2 + \frac{\sigma^2}{n}\right)$

This, together with the prior for  $\eta$  will give the posterior density  $\pi(\eta | \underline{x}) = \pi(\eta | \bar{\underline{x}})$ .

Then the Bayes estimate of  $\mu_j$  is

$$E(\mu_j | \tilde{x}) = E(\mu_j | \bar{x}) = \int E(\mu_j | \bar{x}, \eta) \pi(\eta | \bar{x}) d\eta$$

where,  $E(\mu_j | \bar{x}, \eta) = E(\mu_j | \bar{x}_j, \eta) = (1-B)\bar{x}_j + B\eta_1$  and

$$B = \frac{\sigma^2/n}{\eta_2 + \sigma^2/n}$$

### Parametric Empirical Bayes Approach (PEB):-

Recall that given  $\eta$ ,  $\bar{x}_j$ 's are independent and

$$\bar{x}_j | \eta \sim f(\bar{x}_j | \eta) = N\left(\eta_1, \eta_2 + \frac{\sigma^2}{n}\right) \text{ density.}$$

In this approach, one takes an intermediate position between a fully Bayes and a fully frequentist approach by treating the likelihood as given by  $f(\bar{x}_j | \eta)$  obtained by integrating out the  $\mu_j$ 's. The  $\mu_j$ 's are treated as random variables as in Bayesian analysis whereas  $\eta_1, \eta_2$  are treated as (fixed) unknown parameters as in frequentist analysis. The PEB approach differs from the fully Bayesian (FB) approach in that no prior is assigned to  $\eta$ , and  $\eta$  is estimated by MLE or by a suitable unbiased estimate.

Here, given  $\eta$ ,  $\bar{x}_1, \dots, \bar{x}_p$  are iid  $N\left(\eta_1, \eta_2 + \frac{\sigma^2}{n}\right)$

The best unbiased estimate (UMVUES) of  $\eta_1$  and  $B$  are given

by  $\hat{\eta}_1 = \frac{1}{p} \sum_{j=1}^p \bar{x}_j = \bar{\bar{x}}$ , say and  $\hat{B} = \frac{(p-3)\frac{\sigma^2}{n}}{S}$

If  $Y_1, Y_2, \dots, Y_p$  are iid  $N(\mu, \tau^2)$ ,  $\left[ \because S = \sum_{j=1}^p (\bar{x}_j - \bar{\bar{x}})^2 \right]$   
 UMVUE of  $\mu$  is  $\bar{Y} = \frac{1}{p} \sum Y_j$  and UMVUE of  $\frac{1}{\tau^2}$  is

$$(p-3) / \sum_{j=1}^p (Y_j - \bar{Y})^2$$

Then PEB estimate of  $\mu_j$  is  $\hat{\mu}_j = (1-\hat{B})\bar{x}_j + \hat{B}\bar{\bar{x}}$   
 [a shrinkage estimate shrinking towards  $\bar{\bar{x}}$ ]

If we assume  $\eta_1 = 0$ , i.e., we assume  $\mu_1, \dots, \mu_p$  are iid  $N(0, \eta_2)$  then given  $\eta, \bar{X}_1, \dots, \bar{X}_p$  are iid  $N(0, \eta_2 + \frac{\sigma^2}{n})$ .

Then UMVUE of  $B$  is  $\hat{B} = \frac{(p-2)\frac{\sigma^2}{n}}{\sum_j \bar{X}_j^2}$  and the PEB estimate

of  $\mu_j$  is  $\hat{\mu}_j = \left(1 - \frac{(p-2)\frac{\sigma^2}{n}}{\sum_j \bar{X}_j^2}\right) \bar{X}_j$ , which is the well-known James - Stein estimate.

Remark:- The classical estimate of  $\mu_j$  is  $\bar{X}_j$ , which is based only on the  $j^{\text{th}}$  sample. The Bayes estimate (HB or PEB) of  $\mu_j$ , on the other hand, depends on  $\bar{X}_j$  and also on  $(\bar{X}_1, \dots, \bar{X}_p)$ , the full sufficient statistic because the posterior distribution of  $\eta$  on the estimate of  $\eta$  depends on all the  $\bar{X}_j$ 's. Thus the Bayes estimate learns from the full sufficient statistic justifying simultaneous estimation of all the  $\mu_j$ 's. This learning process is sometimes referred to as borrowing strength. If the modelling of  $\mu_j$ 's is realistic, we would expect the Bayes estimates to perform better than  $\bar{X}_j$ 's. That is what strikingly new in the case of large  $p$ , and follows from the exchangeability of the  $\mu_j$ 's.

## Assignment I

1. Consider the problem of estimation of a real parameter  $\theta$  with the loss function

$$L(\theta, a) = \begin{cases} k_0(\theta - a) & \text{if } \theta - a \geq 0 \\ k_1(a - \theta) & \text{if } \theta - a < 0 \end{cases}$$

s.t. the Bayes estimate is given by the quantile of order  $\frac{k_0}{k_0+k_1}$  of the posterior distr.

Sol.  $L(\theta, a) = k_0(\theta - a) - (k_0+k_1)(\theta - a)I(\theta < a)$  — ①  
 $= k_0(\theta - a) - (k_0+k_1)(\theta - a)I(\theta \leq a)$  — ②

Let  $a_0 \in \mathbb{R}$  s.t.  $P(\theta \leq a_0 | \underline{X}) \geq \frac{k_0}{k_0+k_1}$   
 $P(\theta \geq a_0 | \underline{X}) \geq \frac{k_1}{k_0+k_1}$ .

We are to show  $E_{\theta} [L(\theta, a) | \underline{X}]$  is min. when  $a = a_0$ .

Assume  $a < a_0$ ,

$$L(\theta, a) - L(\theta, a_0)$$

$$= k_0(a_0 - a) - (k_0+k_1)(\theta - a)I(\theta < a) + (k_0+k_1)(\theta - a_0)I(\theta < a_0)$$

[Using ①]

$$= k_0(a_0 - a) - (k_0+k_1)(\theta - a)I(\theta < a_0) + (k_0+k_1)(\theta - a)I(a \leq \theta < a_0)$$

$$+ (k_0+k_1)(\theta - a_0)I(\theta < a_0).$$

$$= k_0(a_0 - a) + (k_0+k_1)(a - a_0)I(\theta < a_0) + \underbrace{(k_0+k_1)(\theta - a)I(a < \theta < a_0)}_{\geq 0}$$

$$\geq (a_0 - a)(k_0+k_1) \left[ \frac{k_0}{k_0+k_1} - I(\theta < a_0) \right] \geq 0$$

$$E [L(\theta, a) - L(\theta, a_0) | \underline{X}]$$

$$\geq (a_0 - a)(k_0+k_1) \left[ \frac{k_0}{k_0+k_1} - P(\theta < a_0 | \underline{X}) \right] > 0$$

Assume  $a > a_0$

$$L(\theta, a) - L(\theta, a_0)$$

$$= k_0(a_0 - a) - (k_0 + k_1)(\theta - a)I(\theta \leq a) + (k_0 + k_1)(\theta - a_0)I(\theta \leq a_0) \quad [\text{Using } \textcircled{2}]$$

$$= k_0(a_0 - a) - (k_0 + k_1)(\theta - a)I(\theta \leq a_0) - (k_0 + k_1)(\theta - a)I(a_0 \leq \theta \leq a) + (k_0 + k_1)(\theta - a_0)I(\theta \leq a_0)$$

$$= (k_0 + k_1)(a - a_0) \left[ I(\theta \leq a_0) - \frac{k_0}{k_0 + k_1} \right] + (k_0 + k_1)(\theta - a)I(a_0 < \theta \leq a)$$

$$\geq (k_0 + k_1)(a - a_0) \left[ I(\theta \leq a_0) - \frac{k_0}{k_0 + k_1} \right] \geq 0$$

$$\therefore E \left[ L(\theta, a) - L(\theta, a_0) \mid \tilde{X} \right]$$

$$\geq (k_0 + k_1)(a - a_0) \left[ P(\theta \leq a_0 \mid \tilde{X}) - \frac{k_0}{k_0 + k_1} \right]$$

$\geq 0$ .

$$\therefore a_0 = \arg \min_a E \left( L(\theta, a) \mid \tilde{X} \right).$$

$\square$ . S.T. the result on asymptotic normality of the posterior distn. of  $\sqrt{n}(\theta - \hat{\theta}_n)$ , proved in the class, implies consistency of the posterior distn. of  $\theta$  at  $\theta_0$ .

Sol. We know that with  $[P_{\theta_0}]$  probability 1

$$\sqrt{n}(\theta - \hat{\theta}_n) \mid \tilde{X}_n \xrightarrow{d} N(0, I^{-1}(\theta_0))$$

$\hat{\theta}_n$  is strongly consistent soln of  $L_n'(\theta) = 0$ .

$\hat{\theta}_n \rightarrow \theta_0$  with  $[P_{\theta_0}]$  prob. 1.

So,  $\hat{\theta}_n \mid \tilde{X}_n \rightarrow \theta_0$  a.s.  $[P_{\theta_0}]$

$$\sqrt{n}(\theta - \hat{\theta}_n) \mid \tilde{X}_n \xrightarrow{d} N(0, I^{-1}(\theta_0)) \text{ a.s. } [P_{\theta_0}]$$

$$\Rightarrow (\theta - \hat{\theta}_n) \mid \tilde{X}_n \xrightarrow{P} 0 \text{ a.s. } [P_{\theta_0}]$$

$$\Rightarrow (\hat{\theta}_n - \theta_0) \mid \tilde{X}_n \xrightarrow{P} 0 \text{ a.s. } [P_{\theta_0}]$$

$$\Rightarrow (\theta - \theta_0) \mid \tilde{X}_n \xrightarrow{P} 0 \text{ a.s. } [P_{\theta_0}]$$

$$\Rightarrow \theta \mid \tilde{X}_n \xrightarrow{P} \theta_0 \text{ a.s. } [P_{\theta_0}]$$

$\Rightarrow$  For any  $U$  open set containing  $\theta_0$   $P[\theta \in U \mid \tilde{X}_n] \rightarrow 1$  a.s.  $[P_{\theta_0}]$

$\Rightarrow \theta \mid \tilde{X}_n$  is consistent at  $\theta_0$ .

5. Show that condition (A4), used to prove asymptotic normality of posterior distn, holds when  $X_1, \dots, X_n$  are i.i.d.  $N(\theta, 1)$ .

Sol. A4 condition: for any  $\delta > 0$ ,  
 $P_{\theta_0} \left[ \text{For some } \epsilon > 0, \sup_{|\theta - \theta_0| > \delta} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)] < -\epsilon \forall n \right]$

for  $X_1, X_2, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} N(\theta, 1) = 1$ .

$$\begin{aligned} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)] &= -\frac{1}{n} \left[ \sum_{i=1}^n [X_i - \theta]^2 - [X_i - \theta_0]^2 \right] \\ &= -\frac{1}{n} \sum_{i=1}^n (2X_i - \theta - \theta_0)(\theta_0 - \theta) \\ &= -(\theta_0 - \theta)(2\bar{X}_n - \theta - \theta_0) \\ &\xrightarrow{\text{a.s.}} -(\theta_0 - \theta)(2\theta_0 - \theta - \theta_0) \\ &= -(\theta_0 - \theta)^2 \text{ under } [P_{\theta_0}]. \end{aligned}$$

for  $\delta > 0$ , for  $\omega \in \Omega$  s.t.

$$\frac{1}{n} [\ln(\theta) - \ln(\theta_0)](\omega) \rightarrow -(\theta - \theta_0)^2$$

then  $\sup_{|\theta - \theta_0| > \delta} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)](\omega) \xrightarrow{\text{a.s.}} \sup_{|\theta - \theta_0| > \delta} -(\theta - \theta_0)^2 = -\delta^2$ .

choose  $\epsilon = \frac{\delta^2}{2} > 0$  then  $\forall n$

$$\frac{1}{n} [\ln(\theta) - \ln(\theta_0)](\omega) < -\epsilon$$

As,  $P_{\theta_0} \left[ \omega \in \Omega : \frac{1}{n} [\ln(\theta) - \ln(\theta_0)](\omega) \rightarrow -(\theta - \theta_0)^2 \right] = 1$

So,  $P_{\theta_0} \left[ \exists \epsilon > 0, \sup_{|\theta - \theta_0| > \delta} \frac{1}{n} [\ln(\theta) - \ln(\theta_0)] < -\epsilon \forall n \right] = 1$ .



10. (a) Let  $X_1, \dots, X_n$  be iid  $\sim N(0, 1)$ . Consider the problem of testing  $H_0: \theta \leq \theta_0$  vs.  $H_1: \theta > \theta_0$ . A classical test rejects  $H_0$  if  $T = \sqrt{n}(\bar{X} - \theta_0)$  is large. Let  $t$  be the observed value of  $T$ . Find the p-value (in terms of  $t$ ).

Consider now the uniform prior  $\pi(\theta) \equiv 1$ . Find the posterior prob. of  $H_0$  (in terms of  $t$ ) and compare it with the p-value.

Soln:-  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N(0, 1)$   
 $\bar{X}_n \sim N(0, 1/n)$  under  $[P_0]$

$\sqrt{n}(\bar{X}_n - \theta) \sim N(0, 1)$  under  $[P_0]$

$$P_\theta [\sqrt{n}(\bar{X} - \theta_0) > t] = P_\theta [\sqrt{n}(\bar{X} - \theta) - \sqrt{n}(\theta_0 - \theta) > t]$$

$$= P_\theta [\sqrt{n}(\bar{X} - \theta) > t + \sqrt{n}(\theta_0 - \theta)]$$

$$= 1 - \Phi(t + \sqrt{n}(\theta_0 - \theta)).$$

$$p = \sup_{\theta \leq \theta_0} (1 - \Phi(t + \sqrt{n}(\theta_0 - \theta)))$$

$$= 1 - \sup_{\theta \leq \theta_0} \Phi(t + \sqrt{n}(\theta_0 - \theta))$$

$$= 1 - \Phi(t).$$

2nd part

Alt:

$$P_{\theta|X}(\theta \leq \theta_0)$$

$$= P\left(N\left(\bar{x}, \frac{\sigma^2}{n}\right) \leq \theta_0\right)$$

$$= P\left(N(0, 1) \leq \frac{\sqrt{n}(\theta_0 - \bar{x})}{\sigma}\right)$$

$$= P(N(0, 1) \leq -t)$$

$$= \Phi(-t)$$

Consider  $\pi(\theta) \equiv 1$ .

$$f(x_1, \dots, x_n | \theta) \propto \exp\left[-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2\right]$$

$$\propto \exp\left[-\frac{n\theta^2}{2} + n\theta\bar{x}_n\right]$$

$$\propto \exp\left[-\frac{n}{2}(\theta - \bar{x}_n)^2\right]$$

$$\pi(\theta | x_1, \dots, x_n) \propto f(x_1, \dots, x_n | \theta) \pi(\theta)$$

$$\propto \exp\left[-\frac{n}{2}(\theta - \bar{x}_n)^2\right]$$

So,  $\theta | x_1, \dots, x_n \sim N(\bar{x}_n, \frac{1}{n})$ .

$$P(H_0 | x_1, \dots, x_n) = P(\theta \leq \theta_0 | x_1, \dots, x_n)$$

$$= P(\sqrt{n}(\theta - \bar{x}_n) \leq \sqrt{n}(\theta_0 - \bar{x}_n) | x_1, \dots, x_n)$$

$$= \Phi(\sqrt{n}(\theta - \bar{x}_n)) = \Phi(-t)$$

$$= 1 - \Phi(t).$$

(11) Let the sample space be  $\{1, 2, \dots, k\}$  and  $P = (p_1, \dots, p_k)$  be a random probab. distr. on this sample space. Let  $X_1, X_2, \dots, X_n$  be iid  $\sim P$ , and  $P \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ ,  $\alpha_i > 0 \forall i$ . S.T. any subset  $A$  of the sample space, the posterior mean of  $P(A)$  is a weighted average of its prior mean and  $P_n(A)$  where,  $P_n$  denotes the empirical distr. of  $X_1, X_2, \dots, X_n$ .

Sol. If  $P \sim D(\alpha_1, \dots, \alpha_k)$ ;  $\alpha_i > 0 \forall i$   
and  $A \subseteq \mathcal{X} = \{1, 2, \dots, k\}$

then  $(p_1, \dots, p_k) \sim D(\alpha_1, \dots, \alpha_k)$

$$P(A) = \sum_{i=1}^k p_i \mathbb{I}(i \in A) \sim \text{Beta}(\alpha(A), \alpha(\mathcal{X}) - \alpha(A)).$$

where  $\alpha(A) = \sum_{i=1}^k \alpha_i \mathbb{I}(i \in A)$  [ follows from the properties of Dirichlet distr. ]

$$\text{So, } E(P(A)) = \frac{\alpha(A)}{\alpha(\mathcal{X})} = \bar{\alpha}(A).$$

Let  $X_1, \dots, X_n \sim P$   
then  $f(X_1, \dots, X_n | P) = p_1^{n_1} \dots p_k^{n_k}$  where  $n_j = \sum_{i=1}^n \mathbb{I}(X_j = i)$

$$P \sim D(\alpha_1, \dots, \alpha_k)$$

$$\Rightarrow \pi(P) \propto p_1^{\alpha_1 - 1} \dots p_k^{\alpha_k - 1}$$

$$\pi(P | X_1, \dots, X_n) \propto \pi(P) f(X_1, \dots, X_n | P)$$

$$= p_1^{n_1 + \alpha_1 - 1} \dots p_k^{n_k + \alpha_k - 1}$$

$$\text{So, } P | X_1, \dots, X_n \sim D(\alpha_1 + n_1, \dots, \alpha_k + n_k)$$

Define  $\lambda: P(\mathcal{X}) \rightarrow \mathbb{R}^+$

Set  $\lambda(\{i\}) = n_i$  for  $i \in \mathcal{X}$ .

$$\text{Then } E(P(A) | X_1, \dots, X_n) = \frac{\alpha(\bar{A}) + \lambda(A)}{\alpha(\mathcal{X}) + n}$$

$$= \frac{\alpha(\bar{\mathcal{X}})}{\alpha(\mathcal{X}) + n} \cdot \frac{\alpha(A)}{\alpha(\mathcal{X})} + \frac{n}{\alpha(\mathcal{X}) + n} \cdot \frac{\lambda(\bar{A})}{n}$$

$\uparrow$   $E(P(A))$ 
 $\uparrow$   $P_n(A)$

⑧ Let  $X_1, X_2$  be a n.s. from a  $U(\theta - 1/2, \theta + 1/2)$  distn.,  $-\infty < \theta < \infty$ . Consider the noninformative prior  $\pi(\theta) = 1$ . What will be your 95% credible interval for  $\theta$ ?

Sol.  $\pi(\theta) = 1$ .

$$f(x_1, x_2 | \theta) = I(\theta - 1/2 < x_{(1)} \leq x_{(2)} < \theta + 1/2)$$

$$\pi(\theta | x_1, x_2) \propto I(\theta - 1/2 < x_{(1)} \leq x_{(2)} < \theta + 1/2)$$

$$\Rightarrow \theta < x_{(1)} + 1/2 \quad \text{and} \quad \theta > x_{(2)} - 1/2.$$

$$\text{Also, } x_{(2)} - 1/2 < \theta < x_{(1)} + 1/2.$$

$$\therefore \pi(\theta | x_1, x_2) = I(x_{(2)} - 1/2 < \theta < x_{(1)} + 1/2)$$

$$= \frac{1}{x_{(1)} - x_{(2)}} = 1 - |x_1 - x_2|$$

95% credible set is any

$$C \subseteq (x_{(2)} - 1/2, x_{(1)} + 1/2)$$

$$\text{s.t. } P(\theta \in C | x_1, x_2) = 0.95$$

Let  $C = (\bar{x} - k, \bar{x} + k)$  then

$$P(\theta \in C | x_1, x_2) = \frac{2k}{1 - |x_1 - x_2|} = 0.95$$

$$\Rightarrow k = \frac{0.95(1 - |x_1 - x_2|)}{2}$$

$$\text{So, } C = \left( \bar{x} \mp \frac{0.95(1 - |x_1 - x_2|)}{2} \right).$$