# Paradoxes in Classical Statistics

**Example 1.** To estimate $\mu$ in $N(\mu, \sigma^2)$, toss a fair coin. Have a sample size $n = 2$ if it is a head and take $n = 1000$ if it is a tail. An unbiased estimate of $\mu$ is $\bar{X}_n = \sum_{i=1}^{n} X_i / n$ with variance

$$\frac{1}{2} \left\{ \frac{\sigma^2}{2} + \frac{\sigma^2}{1000} \right\} \approx \frac{\sigma^2}{4}.$$

Suppose it was a tail. Would you believe $\frac{\sigma^2}{4}$ is a measure of accuracy of the estimate?

**Example 2.** Let $X_1, X_2$ be i.i.d. $U[\theta - \frac{1}{2}, \theta + \frac{1}{2}]$. Let $\bar{X} \pm C$ be a 95% confidence interval, $C > 0$ being suitably chosen. Suppose $X_1 = 2$ and $X_2 = 1$. Then we know for sure $\theta = (X_1 + X_2)/2$ and hence $\theta \in (\bar{X} - C, \bar{X} + C)$ Should we still claim we have only 95% confidence that the confidence interval covers $\theta$ ?

## Bayesian Inference

To make inference about $\theta$ is to learn about the unknown $\theta$ from data $X$, i.e, based on the data, explore which values of $\theta$ are probable, what might be estimates of different components of $\theta$ and the extent of uncertainty associated with such estimates.

In addition to having a model $f(x|\theta)$ and a likelihood function, the Bayesian needs a distribution to ~~having~~ a for $\theta$.

The distribution is called a prior distribution or simply a prior because it quantifies her uncertainty about $\theta$ prior to seeing data.

Now the conditional probability density of $\theta$ given $X = x$ by Bayes formula

$$\pi\left(\theta \mid x\right) = \frac{\pi(\theta)\, f\,(x \mid \theta)}{\int_{\Theta} \pi(\theta')\, f(x \mid \theta')\, d\theta'} \qquad \cdots \cdots (1)$$

where $\pi(\theta)$ is the prior density function and $f(x \mid \theta)$ is the density of $X$, which is the conditional density of $X$ given $\theta$.

The symbol $\theta$ now represents both a random variable and its value. When $\theta$ is discrete, the integral in the denominator of (1) is replaced by a sum.

The Conditional density $\pi(\theta \mid x)$ of $\theta$ given $X = x$ is called the posterior density, a quantification of our uncertainty about $\theta$ in the light of data. The transition from $\pi(\theta)$ to $\pi(\theta \mid x)$ is what we learnt from the data.

A Bayesian can simply report her posterior distribution, or she could report the posterior mean

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

and the posterior variance

$$Var(\theta | x) = \int_{-\infty}^{\infty} (\theta - E(\theta | x))^2 \, \pi(\theta | x) \, d\theta.$$

if she want to estimate $\theta$,

or in the case of testing she would report the posterior odds of the relevant hypotheses.

<u>Example 1</u> Consider the problem of inference about $\mu$ for normally distributed data $N(\mu, \sigma^2)$.

The data consist of i.i.d. observation $X_1, X_2, \cdots, X_n$ from this distribution. We assume that $\sigma^2$ is known.

A mathematically convenient and reasonably flexible prior distribution for $\mu$ is a normal distribution with suitable prior mean $\eta$ and variance $\tau^2$. The prior variance, $\tau^2$ is a measure of strength of our belief in the prior mean, in the sense that the larger the value of $\tau^2$, the less sure we are about our prior guess about $\eta$. Jeffreys has suggested we can calibrate $\tau^2$ by comparing with $\sigma^2$.

For example setting $\tau^2 = \sigma^2/m$ would amount to saying information about $\eta$ is about as strong as the information in $m$ observations in data.

To illustrate, we take, $n = 10$, $m = 1$, so the posterior distribution is normal with posterior mean

$$E(\mu \mid x) = \left(\frac{1}{\tau^2}\eta + \frac{n}{\sigma^2}\bar{x}\right)\Big/\left(\frac{1}{\tau^2} + \frac{n}{\sigma^2}\right)$$

$$= (\eta + 10\,\bar{x})/11$$

and posterior variance

$$\left(\frac{\sigma^2}{n}\,\tau^2\right)\Big/\left(\frac{\sigma^2}{n} + \tau^2\right) = \frac{\sigma^2}{11}$$

i.e., in the light of data, $\mu$ shifts from prior guess $\eta$ towards a weighted average of the prior guess about $\mu$ and $\bar{x}$, while the variability reduces from $\sigma^2$ to $\sigma^2/11$.

If the prior information is small, implying large $\tau^2$ or there are lots of data i.e., $n$ is large, the posterior mean is close to the MLE $\bar{x}$.

**Example 2** Consider an urn with $Np$ red and $N(1-p)$ black balls, $p$ is unknown but $N$ is a known large number. Balls are drawn at random one by one with replacement, selection is stopped after $n$ draws. For $i = 1, 2, \cdots, n$ let

$$X_i = \begin{cases} 1 & \text{if the } i\text{th ball drawn is red} \\ 0 & \cdot \text{otherwise}. \end{cases}$$

Then $X_i$'s are i.i.d. $B(1, p)$,

Let $p$ have a prior distribution $\pi(p)$.

We will consider a family of priors for $p$ that simplifies the calculation of posterior. 

Let $$\pi(p) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot p^{\alpha-1}(1-p)^{\beta-1}, \quad \cdots \cdots (2)$$

$$0 \le p \le 1 \; ; \; \alpha > 0, \; \beta > 0$$

The prior mean and variance are

$$\frac{\alpha}{\alpha+\beta} \quad \text{and} \quad \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}, \quad \text{respectively}.$$

The posterior density is

$$\pi(p \mid X = x) = C(x) \, p^{\alpha+r-1}(1-p)^{\beta+(n-r)-1}$$

where $r = \sum_{i=1}^{n} x_i$ = number of red balls in $n$ draws and $(C(x))^{-1}$ is the denominator in Bayes formula.

Now the posterior is also a Beta density with $\alpha + r$ in place of $\alpha$ and $\beta + (n - r)$ in place of $\beta$ in (2) and

$$C(x) = \frac{\Gamma(\alpha + \beta + n)}{\Gamma(\alpha + r)\,\Gamma(\beta + n - r)}$$

The posterior mean and variance are

$$E(p\,|\,x) = (\alpha + r)\big/(\alpha + \beta + n)$$

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

Again if $n$ is large then the posterior mean is approximately equal to the MLE, $\hat{p} = r/n$ and the posterior variance is quite small,

We can interpret this as an illustration of a fact mentioned before when we have lots of data, the data tend to wash away the influence of the prior.

The posterior mean can be rewritten as a weighted average of the prior mean and MLE

$$\frac{(\alpha + \beta)}{(\alpha + \beta + n)} \cdot \frac{\alpha}{(\alpha + \beta)} + \frac{n}{(\alpha + \beta + n)} \cdot \frac{r}{n}$$

the relative importance of the prior and the data being measured by $(\alpha + \beta)$ and $n$.

∴ Suppose we want to predict the probability of getting a red ball in a new $(n+1)$-st draw given the above data. It would be natural to use $E(p/\lambda)$.

Some commonly used priors and the corresponding value of $E(P/X_1, X_2, \cdots, X_n)$

If $\alpha = \beta = 1$, the uniform prior, with posterior mean $\left(\sum_{i=1}^{n} x_i + 1\right)/(n+2)$.

If $\alpha = \beta = \frac{1}{2}$, we have the Jeffreys prior with posterior mean $\left(\sum_{i}^{n} x_i + \frac{1}{2}\right)/(n+1)$

<u>Example 3</u>   Suppose we think of the problem as a representation of production of defective and non-defective items in a factory producing switches, we would take red to mean defective and black to mean a good switch.

In this context, there would be some prior information available from the engineers. They may be able to pinpoint the likely value of $p$, which may be set equal to the prior mean $\alpha/(\alpha+\beta)$. If one has Knowledge of prior variability also, one would have two equations from which to determine $\alpha$ and $\beta$.

In this particular context, the Jeffreys prior with a lot It can be shown that the uniform, Jeffreys prior produce a posterior mean that is very closed to the MLE even for small $n$.

Also they make better sense than the MLE in the extreme case when $\hat{p} = 0$. In most contexts the estimate $\hat{p} = 0$ is absurd, the objective Bayes estimates move it a little towards $p = \frac{1}{2}$. Such a movement is called a shrinkage.

### Example 4

Let $X_1, \ldots X_n$ k.i.i.d $N(\mu, \sigma^2)$ and $\sigma^2$ is known.

We want to test $H_0 : \mu \leq \mu_0$ vs $H_1 : \mu > \mu_0$.

Let $\pi(\mu)$ be the prior. First calculate the posterior density $\pi(\mu \mid x)$. Then

$$P\{H_0 \mid x\} = \int_{-\infty}^{\mu_0} \pi(\mu \mid x) \, d\mu$$

and $\int_{\mu_0}^{\infty} \pi(\mu \mid x) \, d\mu = 1 - P\{H_0 \mid x\} = P(H_1 \mid x)$

One may simply report these numbers or choose one of the two hypotheses if one of the two probabilities is substantially bigger.

## Improper Priors

For point and interval estimates and to some extent in testing improper priors are used.

An improper prior density $\pi(\theta)$ is non-negative for all $\theta$ but $\int_\Theta \pi(\theta)\, d\theta = \infty$.

Such an improper prior can be used in the Bayes formula for calculating the posterior, provided the denominator is finite for all $x$ (or all but a set of $x$ with zero probability) and for all $\theta$, i.e.,

$$\int_\Theta \pi(\theta)\, f(x/\theta)\, d\theta < \infty.$$

Then the posterior density $\pi(\theta / x = x)$ is a proper density function and can be used at least in inference problems or the posterior density function decision problem where we define and minimize

$$\psi(x,a) = E\left(L(\theta,a)/x\right) = \int_\Theta L(\theta,a)\, \pi(\theta/x)\, d\theta$$

which is the posterior risk.

However, for improper priors usually $R(\pi, \delta)$ is not used.

where $R(\pi, \delta) = \int_\Theta R(\theta, \delta)\, \pi(\theta)\, d\theta$

and $R(\theta, \delta) = E_\theta \left(L(\theta, \delta(x))\right)$

The most common improper priors are

$$\pi_1(\mu) = C, \qquad -\infty < \mu < \infty,$$

$$\pi_2(\sigma) = \frac{1}{\sigma}, \qquad 0 < \sigma < \infty$$

for location and scale parameters $\mu$ and $\sigma$ resp.

Both the improper priors may be interpreted as a port of limit of the proper priors.

$$\pi_{1,L}(\mu) = \begin{cases} \frac{1}{2L} & \text{if } -L < \mu < L \\ 0 & \text{otherwise} \end{cases}$$

and

$$\pi_{2;L}(\sigma) = \begin{cases} A/\sigma & \text{if } 0 < \frac{1}{L} < \sigma < L \\ 0 & \text{otherwise}; \end{cases}$$

where $A = \frac{1}{(2 \log L)}$, in the sense that the posteriors for $\pi_1$ and $\pi_2$ may be obtained by making $L \to \infty$ in $\pi_{i,L}(\theta \mid x)$

## Point Estimates

For a real value θ, standard Bayes estimates are the posterior mean or the posterior median.

The posterior man is the Bayes estimate corresponding with squared error loss and the posterior median is Bayes estimate for absolute deviation loss.

Along with the posterior mean one reports the posterior variance of θ. If one choose to work with the posterior median, one should report atleast first and third posterior quartiles.

If p the posterior is unimodal then the posterior mode is another choice. It is similar to the MLE of classical statistics. Indeed if the prior is uniform, both are identical. Along with the posterior mode one can report a suitable highest posterior density credible interval as a measure of posterior variability.

If the parameter is a vector, common choices for reporting are the posterior mean vector and the posterior dispersion matrix.

# Testing

We want to test $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$. If $\Theta_0$ and $\Theta_1$ are of the same dimension as for one-sided null and alternative hypothesis, it is convenient and easy to choose a prior distribution that assigns positive prior probability to $\Theta_0$ and $\Theta_1$. One then calculates the posterior probabilities $P\{\Theta_i | x\}$ as well as the posterior odds ratio (or simply posterior odds), namely,

$$P\{\Theta_0 | x\} \Big/ P\{\Theta_1 | x\}$$

One would then find a threshold to decide what constitutes evidence against $H_0$. The Bayes rule for 0-1 loss is to choose the hypothesis with higher posterior probability.

Let $\pi_0$ and $1 - \pi_0$ be the prior probabilites of $\Theta_0$ and $\Theta_1$. Let $g_i(\theta)$ be the prior p.d.f. of $\theta$ under $\Theta_i$, so that $\int_{\Theta_i} g_i(\theta) \, d\theta = 1$

The prior is then.

$$\pi(\theta) = \pi_0 \, g_0(\theta) \, I\{\theta \in \Theta_0\} + (1 - \pi_0) \, g_1(\theta) \, I\{\theta \in \Theta_1\}.$$

We do not require any longer that $\Theta_0$ and $\Theta_1$ are the same dimension. So sharp null hypotheses are also covered. We can now proceed as before and report posterior probabilities or posterior odds.

To compute these posterior quantities, note that the marginal density of $x$ under the prior $\pi$ can be expressed as

$$m_\pi(x) = \int_\Theta f(x|\theta)\,\pi(\theta)\,d\theta$$

$$= \pi_0 \int_{\Theta_0} f(x|\theta)\,g_0(\theta)\,d\theta + (1-\pi_0)\int_{\Theta_1} f(x|\theta)\,g_1(\theta)\,d\theta$$

and hence the posterior density of $\theta$ given the data $X = x$ as

$$\pi(\theta|x) = \frac{f(x|\theta)\,\pi(\theta)}{m_\pi(x)}$$

$$\therefore \quad \pi(\theta|x) = \begin{cases} \pi_0\,f(x|\theta)\,g_0(\theta)/m_\pi(x) & \text{if } \theta \in \Theta_0 \\[2ex] (1-\pi_0)\,f(x|\theta)\,g_1(\theta)/m_\pi(x) & \text{if } \theta \in \Theta_1. \end{cases}$$

It follows that

$$P^\pi(H_0|x) = P^\pi(\Theta_0|x) = \frac{\pi_0}{m_\pi(x)}\int_{\Theta_0} f(x|\theta)\,g_0(\theta)\,d\theta$$

and

$$P^\pi(H_1|x) = P^\pi(\Theta_1|x) = \frac{(1-\pi_0)}{m_\pi(x)}\int_{\Theta_1} f(x|\theta)\,g_1(\theta)\,d\theta$$

Now the Bayes factor of $H_0$ relative to $H_1$ is defined as

$$BF_{01} = \frac{\displaystyle\int_{\Theta_0} f(x|\theta)\,g_0(\theta)\,d\theta}{\displaystyle\int_{\Theta_1} f(x|\theta)\,g_1(\theta)\,d\theta}$$

which does not depends on $\pi_0$.

clearly $BF_{10} = \dfrac{1}{BF_{01}}$

The posterior odds ratio of $H_0$ relative to $H_1$ is

$$\left( \frac{\pi_0}{1 - \pi_0} \right) BF_{01},$$

which is same as $BF_{01}$ if $\pi_0 = \frac{1}{2}$.

The smaller the value of $BF_{01}$, the stronger the evidence against $H_0$.

Example    Consider a blood test conducted for determining the sugar level of a person with diabetes two hours after he had his breakfast. It is of interest to see if his medication has controlled his blood sugar level. Assume that the test result $X$ is $N(\theta, 100)$, where $\theta$ is the true level. In the appropriate population (diabetic but under this treatment), $\theta$ is distributed according to a $N(100, 900)$. Then marginally $X$ is $N(100, 1000)$, and the posterior distribution of $\theta$ given $X = x$ is normal with mean $\frac{900}{1000} x + \frac{100}{1000} 1000$

$$= 0.9x + 10 \quad \text{and} \quad \text{variance} = \frac{100 \times 900}{1000} = 90$$

Suppose we want to test $H_0 : \theta \leq 130$ vs $H_1 : \theta > 130$. If the blood test shows a sugar level of 130, what can be concluded?

Note that, given this test result, the true-mean blood sugar level, $\theta$, may be assumed to be $N(127, 90)$. Consequently we obtain,

$$P(\theta \leq 130 | x = 130) = \Phi\left( \frac{130 - 127}{\sqrt{90}} \right) = 0.624.$$

Therefore Posterior odds ratio $= \frac{0.624}{1 - 0.624} = 1.66$.

Bacuse $\pi_0 = P^\pi(\theta \leq 130) = \Phi(1) = 0.8413$,

the prior odds ratio is $\dfrac{\Phi(1)}{1 - \Phi(1)} = 5.3$ and thus

the Bayes factor turns out to be $\dfrac{1.66}{5.3} = 0.313$.

It can be noted here that in one-sided testing situations when a continuous prior $\pi$ can be specified readily for the entire parameter space, there is no need to express it in form of

$$\pi(\theta) = \pi_0 \, g_0(\theta) \, I\{\theta \in \Theta_0\} + (1 - \pi_0) \, g_1(\theta) \, I\{\theta \in \Theta_1\}.$$

However, the problem of testing a point null hypothesis turns out to be quite different.

## Testing of a Point Null Hypothesis

The problem is to test

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

Consider the following example.

In a statistical quality control situation, $\theta$ is the size of a unit and acceptable units are with $\theta \in (\theta_0 - \delta, \theta_0 + \delta)$. Then one would like to test $H_{0\delta} : |\theta - \theta_0| \leq \delta$.

In this problem the length of the interval, $2\delta$, can be explicitly specified. On the other hand, this is not in the case of the following example.

Suppose we want to test the hypothesis,

$H_0$ : Vitamin C has no effect on the common cold.
Clearly this is not meant to be thought of as an
exact point null, surely vitamin C has 'some'
effect, though perhaps a very minuscule effect.
Thus, in reality, this is still the case of an
interval null hypothesis, with a very small unspecified
interval. However, it would be better represented
as a point null hypothesis.

Now consider the hypothesis such as

$H_0$ : Astrology cannot predict the future.
Can perhaps be represented as an exact point null.

Conceptually testing a point null is not a
different problem, but there are complications.
First of all, it is not possible to use a continuous
prior density because any such prior will necessarily
assign prior probability zero to the null hypothesis.
Consequently, the posterior probability of the
null hypothesis will also be zero. Therefore,
a prior probability of $\pi_0 > 0$ needs to be assigned
to the point $\theta_0$ and the remaining probability
of $\pi_1 = 1 - \pi_0$ will be spread over $\{\theta \neq \theta_0\}$ using
a density $g_1$.

(9)

simply take $\theta_0$ to be a point mass at $\theta_0$.

∴ If the point null hypothesis approximates an interval ~~via~~ null hypothesis, $H_0: \theta \in (\theta_0 - \epsilon, \theta_0 + \epsilon)$, then $\pi_0$ is the probability assigned to the interval $(\theta_0 - \epsilon, \theta_0 + \epsilon)$ by a continuous prior.

The complication now is that the prior $\pi$ is of the form $\pi(\theta) = \pi_0 I\{\theta = \theta_0\} + (1 - \pi_0) g_1(\theta) I\{\theta \neq \theta_0\}$ and hence has both discrete and continuous part.

Now the marginal density of $X$ is

$$m(x) = \pi_0 f(x/\theta_0) + (1 - \pi_0) m_1(x)$$

where $$m_1(x) = \int_{\theta \neq \theta_0} f(x/\theta) g_1(\theta) d\theta$$

Therefore, $$\pi(\theta_0/x) = \frac{f(x/\theta_0) \pi_0}{m(x)}$$

$$= \frac{\pi_0 f(x/\theta_0)}{\pi_0 f(x/\theta_0) + (1 - \pi_0) m_1(x)}$$

$$= \left\{ 1 + \frac{1 - \pi_0}{\pi_0} \frac{m_1(x)}{f(x/\theta_0)} \right\}^{-1}.$$

It follows then the posterior odds ratio is given by

$$\frac{\pi(\theta_0/x)}{1 - \pi(\theta_0/x)} = \frac{\pi_0}{(1 - \pi_0)} \frac{f(x/\theta_0)}{m_1(x)},$$

and hence the Bayes factor of $H_0$ relative to $H_1$, which is the ratio of the above posterior odds ratio to the prior odds ratio of $\pi_0/(1 - \pi_0)$ is

$$B = B(x) = BF_{01}(x) = \frac{f(x \mid \theta_0)}{m_1(x)}$$

Thus, $\pi(\theta_0 \mid x) = \left\{ 1 + \frac{1 - \pi_0}{\pi_0} BF_{01}^{-1}(x) \right\}^{-1}$

- **Example**

  Suppose $x \sim B(n, \theta)$ and we want to test
  $H_0: \theta = \theta_0$ vs $H_1: \theta \neq \theta_0$. Under the alternative
  hypothesis, suppose $\theta \sim Beta(\alpha, \beta)$. Then $m_1(x)$

  $$m_1(x) = \binom{n}{x} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha+x)\Gamma(\beta+n-x)}{\Gamma(\alpha+\beta+n)}$$

  so that, $BF_{01}(x) = \dfrac{\binom{n}{x} \theta_0^x (1-\theta_0)^{n-x}}{\dfrac{\binom{n}{x}}{} \Gamma \quad m_1(x)}$

  $$= \frac{\theta_0^x (1-\theta_0)^{n-x}}{\dfrac{\Gamma(\alpha+\beta)}{\Gamma(\cdot)}}$$

  $$= \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} \frac{\Gamma(\alpha+\beta+n)}{\Gamma(\alpha+x)\Gamma(\beta+n-x)} \theta_0^x (1-\theta_0)^{n-x}$$

Hence we obtain,

$$\pi(\theta_0 \mid x) = \left\{ 1 + \frac{1-\pi_0}{\pi_0} BF_{01}^{-1}(x) \right\}^{-1}$$

$$= \left\{ 1 + \frac{1-\pi_0}{\pi_0} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \cdot \frac{\Gamma(\alpha+x)\Gamma(\beta+n-x)}{\Gamma(\alpha+\beta+n)} \cdot \frac{1}{\theta_0^x (1-\theta_0)^{n-x}} \right\}^{-1}$$

## Credible Intervals

For $0 < \alpha < 1$, a $100(1-\alpha)\%$ credible set for $\theta$ is a subset $C \subseteq \Theta$ such that

$$P\{C \mid x = x\} = 1 - \alpha$$

Usually $C$ is taken to be an interval.

Let $\theta$ be a continuous random variable; $\theta^{(1)}, \theta^{(2)}$ be $100\,\alpha_1\%$ and $100(1-\alpha_2)\%$ quantiles with $\alpha_1 + \alpha_2 = \alpha$.

Let $C = [\theta^{(1)}, \theta^{(2)}]$. Then $P(C \mid x = x) = 1 - \alpha$.

Usually $\alpha_1 = \alpha_2 = \alpha/2$ are choosen.

If $\theta$ is discrete, usually it would be difficult to find an interval with exact posterior probability $1 - \alpha$. There the condition is relaxed to $P(C \mid x = x) \geqslant 1 - \alpha$, with the inequality being as closed to an equality as possible.

## HPD (Highbest Posterior Density) interval

Suppose the posterior density for $\theta$ is unimodal. Then the HPD interval for $\theta$ is the interval

$$C = \{\theta : \pi(\theta \mid x = x) \geqslant k\}$$

where $k$ is choosen such that

$$P(C \mid x = x) = 1 - \alpha$$

# Prediction of Future Observation

Suppose the data are $x_1, \cdots, x_n$ where $x_1, \cdots, x_n$ are i.i.d. with density $f(x \mid \theta)$. We want to predict the unobserved $x_{n+1}$ or set up a predictive credible interval for $x_{n+1}$.

Prediction by a single number $t(x_1, \cdots, x_n)$ based on $x_1, \cdots, x_n$ with squared error loss amounts to considering prediction loss

$$E\left\{ (x_{n+1} - t)^2 \mid x \right\} = E\left[ \left\{ (x_{n+1} - E(x_{n+1} \mid x)) - (t - E(x_{n+1} \mid x)) \right\}^2 \mid x \right]$$

$$= E\left\{ (x_{n+1} - E(x_{n+1} \mid x))^2 \mid x \right\} + (t - E(x_{n+1} \mid x))^2$$

which is minimum at $t = E(x_{n+1} \mid x)$.

To calculate the predictor we need to calculate the predictive distribution

$$\pi(x_{n+1} \mid x) = \int_{\Theta} \pi(x_{n+1} \mid x, \theta) \, \pi(\theta \mid x) \, d\theta$$

$$= \int_{\Theta} f(x_{n+1} \mid \theta) \, \pi(\theta \mid x) \, d\theta$$

Let $\mu(\theta) = \int_{-\infty}^{\infty} x \, f(x \mid \theta) \, dx$

Then $E(x_{n+1} \mid x) = E(\mu(\theta) \mid x) = \int_{\Theta} \mu(\theta) \, \pi(\theta \mid x) \, dx$

Now if $f(x \mid \theta)$ is a $N(\mu, \sigma^2)$ density with $\sigma^2$ known

$\mu(\theta) = \mu$ and hence the predictor is

$$\int_{-\infty}^{\infty} \mu \, \pi(\mu \mid x) \, d\mu = \text{posterior mean of } \mu.$$

## Exchangeability

A set of observation $\{x_1, x_2, \ldots, x_n\}$ are so called exchangeable, if their joint distribution function is left unaltered if the arguments are permuted.

Thus if i.e., $P\{X_1 \le x_1, \ldots, X_n \le x_n\} = P\{X_1 \le x_{i_1}, \ldots X_n \le x_{i_n}\}$

for all $n!$ permutations $x_{i_1}, \ldots, x_{i_n}$ of $x_1, \ldots x_n$.

Clearly, if $X_1, X_2, \ldots X_n$ are i.i.d. then they are exchangeable. A simple way of generating exchangeable random variables is to choose an indexing random parameter $\eta$ and have the random variables conditionally i.i.d given $\eta$.

We say $X_i$, $i=1,2,\ldots,n,\ldots$ is a sequence of exchangeable random variables if $\forall n > 1$,

$X_1, X_2, \ldots, X_n$ are exchangeable.

__Theorem__. Suppose $X_i$'s constitute an exchangeable sequence and each $X_i$ takes only values 0 or 1. Then, for some $\pi$,

$$P\{X_1 = x_1, \ldots, X_n = x_n\} = \int_0^1 \eta^{\sum_{i=1}^n x_i} (1-\eta)^{n - \sum_{i=1}^n x_i} d\pi(\eta)$$

$\forall n$, $\forall x_1, \ldots, x_n$ equal to 0 or 1, i.e., given $\eta$, $X_1, X_2, \ldots, X_n$ are conditionally i.i.d. Bernoulli with parameter $\eta$ and $\eta$ has distribution $\pi$.

Exercise 1. a Three prisoners, A, B, and C, are each held in solitary confinement. A knows that two of them will be hanged and one will be set free but he does not know who will go free. Therefore, he reasons that he has 1/3 chance of survival. He asks the guard who will go free, but has no success there, then he comes up with the following question for the guard:

If two of us must die, then I know that either B or C must die and possibly both. Therefore, if you tell me the name of one who is to die, I learn nothing about my fate; futher, because we are kept apart, I cannot reveal it to them. So tell me the name of one of them who is to die.

The guard like this logic and tells A that C will be hanged. A now argues that either he or B will go free, and so now has 1/2 chance of survival. Is this reasoning correct?

2. There are three chambers, one of which has a prize. The master of ceremonies will give the prize to you if you guess the right chamber correctly. You first make a random guess. Then he shows you one chamber which is empty (The chamber you guessed first has not been opened). You have an option to stick to your original guess or switch to the remaining other chamber. what should you do.

**Example 3**

Suppose $X|\mu \sim N(\mu, \sigma^2)$, $\sigma^2$ known and $\mu \sim N(\eta, \tau^2)$, $\eta$ and $\tau^2$ known

(a) Show that the joint density $g(x, \mu)$ of $X$ and $\mu$ can be written as

$$g(x, \mu) = \pi(\mu) f(x|\mu) = \frac{1}{2\pi \sigma \tau} \exp\left\{-\frac{1}{2}\left[\frac{(\mu-\eta)^2}{\tau^2} + \frac{(x-\mu)^2}{\sigma^2}\right]\right\}$$

$$= \frac{1}{\sqrt{2\pi(\tau^2+\sigma^2)}} \exp\left(-\frac{(x-\eta)^2}{2(\tau^2+\sigma^2)}\right) \times \sqrt{\frac{\tau^2+\sigma^2}{2\pi \tau^2 \sigma^2}} \exp\left\{-\frac{\tau^2+\sigma^2}{2\tau^2\sigma^2}\left(\mu - \frac{\tau^2\sigma^2}{\tau^2+\sigma^2}\left(\frac{\eta}{\tau^2}+\frac{x}{\sigma^2}\right)\right)^2\right\}$$

(b) Show that the marginal density $m(x)$ of $X$ is

$$m(x) = \frac{1}{\sqrt{2\pi(\tau^2+\sigma^2)}} \exp\left(-\frac{(x-\eta)^2}{2(\tau^2+\sigma^2)}\right),$$

and the posterior density $\pi(\mu|x)$ of $\mu|x = x$ is

$$\pi(\mu|x) = \sqrt{\frac{\tau^2+\sigma^2}{2\pi \tau^2 \sigma^2}} \exp\left\{-\frac{\tau^2+\sigma^2}{2\tau^2\sigma^2}\left(\mu - \frac{\tau^2\sigma^2}{\tau^2+\sigma^2}\left(\frac{\eta}{\tau^2}+\frac{x}{\sigma^2}\right)\right)^2\right\}.$$

(c) What are the posterior mean and posterior s.d. of $\mu$ given $X = x$?

(d) Instead of a single observation $X$ as above, consider a random sample $X_1, \ldots, X_n$. What is the minimal sufficient statistics and what is the likelihood function for $\mu$ now?

Work out (b) and (c) in this case.

4. Suppose $X_1$ and $X_2$ are i.i.d. having the discrete distribution:

$$X = \begin{cases} \theta - \frac{1}{2} & \text{w. p. } \frac{1}{2} \\ \theta + \frac{1}{2} & \text{w. p. } \frac{1}{2} \end{cases}$$

where $\theta$ is an unknown real number.

(a) Show that the set $C$ given by

$$C = \begin{cases} \{(x_1 + x_2)/2\} & \text{if } x_1 \neq x_2 \\ \{x_1 - 1\} & \text{if } x_1 = x_2 \end{cases}$$

is a 75% confidence set for $\theta$.

(b) Calculate $P\{C \text{ covers } \theta \mid x_1 - x_2\}$.

5. Let $X_1, X_2, \ldots X_n$ be i.i.d $N(\mu, \sigma^2)$, $\sigma^2$ known.
Suppose $\mu$ has the $N(\eta, \tau^2)$ prior distribution with known $\eta$ and $\tau^2$

(a) Construct the $100(1-\alpha)\%$ HPD credible interval for $\mu$.

(b) Construct a $100(1-\alpha)\%$ predictive interval for $X_{n+1}$.

(c) Consider the uniform prior for this problem by letting $\tau^2 \rightarrow \infty$. Work out (a) and (b) in this case.

6. Let $X_1, \ldots X_m$ and $Y_1, \ldots Y_n$ be independent random samples, respectively from $N(\mu_1, \sigma^2)$ and $N(\mu_2, \sigma^2)$ where $\sigma^2$ is known. Construct a $100(1-\alpha)\%$ credible interval for $(\mu_1 - \mu_2)$ assuming a uniform prior on $(\mu_1, \mu_2)$.

## Coherence

There is an alternative way of justifying a Bayesian approach to decision making on the basis of the notion of coherence.

Suppose A stands for a set in the space of $\theta$ and $x$ of the and $A_x = \{\theta : (\theta, x) \in A\}$. Given $x$, the decision maker (DM)'s uncertainity about A is given by $q(x, A_x)$.

An MC (master of ceremonines) chooses a betting system $(A, b)$, where A is as above and b is a bounded real valued function of $x$. The DM accepts the ~~gro~~ gamble with pay-off

$$\psi(\theta, x) = b(x)\left[ I_A(\theta, x) - q(x, A_x)\right].$$

She gets $b(x) \, q(x, A_x)$ or pays $b(x)\left[1 - q(x, A_x)\right]$ depending on whether $\theta$ lies in $A_x$ or not. The expected pay-off is

$$E(\theta) = \int \psi(\theta, x) \, p(dx/\theta)$$

If she accepts k such gambles defined as above by $(A_{(1)}, b_{(1)}), \cdots, (A_{(k)}, b_{(k)})$, then her expected pay off is the sum of the k expected pay-offs. She will face sure loss if

$$\inf_\theta \left( \sum_{i=1}^{k} E_i(\theta)\right) > 0$$

Any rational choice of q should avoid sure loss as defined above. Such a choice is said to be coherent if no finite combination of acceptable bets can lead to sure loss.

# Robustness and Sensitivity

Intuitively, robustness means lack of sensitivity of the decision or inference to assumptions in the analysis that may involve a certain degree of uncertainity. In an inference problem, the assumptions usually involve choice of the model and prior, whereas in a decision problem there is the additional assumption involving the choice of the loss or utility function. An analysis to measure the sensitivity is called sensitivity analysis. Clearly, robustness with respect to all ~~three~~ of the components is desirable. That is to say that reasonable variations from the choice used in the analysis for the model, prior, and loss function do not lead to unreasonable variations in the conclusions arrived at.

Example    This example illustrates why sensitivity to the choice of prior can be an important consideration. Suppose we observe X, which follows Poisson($\theta$) distribution. Further, it is felt a priori that $\theta$ has a continuous distribution with median 2 and upper quartile 4, i.e. $P^{\pi}(\theta \leq 2) = 0.5 = P^{\pi}(\theta \geq 2)$

and $P^{\pi}(\theta \geq 4) = 0.25$.

If these are the only prior inputs available, the following three are candidates for such a prior:

(i) $\pi_1$ : $\theta \sim$ exponential $(a)$ with $a = \dfrac{\log 2}{2}$

(ii) $\pi_2$ : $\log(\theta) \sim N\left(\log(2), \left(\dfrac{(\log 2)}{z_{0.25}}\right)^2\right)$; and

(iii) $\pi_3$ : $\log(\theta) \sim$ Cauchy $(\log(2), \log(2))$.

Then (i) Under $\pi_1$, $\theta|x \sim$ Gamma $(a+1, x+1)$,

so that the posterior mean is $\dfrac{(a+1)}{(x+1)}$.

(ii) under $\pi_2$, if we let $\gamma = \log(\theta)$ and

$$\tau = \dfrac{\log(2)}{z_{0.25}} = \dfrac{\log(2)}{0.675}, \quad \text{we} \quad \text{obtain}$$

$$E^{\pi_2}(\theta|x) = E^{\pi_2}(\exp(\gamma)|x)$$

$$= \dfrac{\displaystyle\int_{-\infty}^{\infty} \exp(-e^{\gamma})\, \exp(\gamma(x+1))\, \exp(-(\gamma-\log 2)^2/2\tau^2)\, d\gamma}{\displaystyle\int_{-\infty}^{\infty} \exp(-e^{\gamma})\, \exp(\gamma x)\, \exp(-(\gamma-\log 2)^2/2\tau^2)\, d\gamma}$$

and (iii) under $\pi_3$, again if let $\gamma = \log\theta$ we get

$$E^{\pi_3}(\theta|x) = E^{\pi_3}(\exp(\gamma)|x)$$

$$= \dfrac{\displaystyle\int_{-\infty}^{\infty} \exp(-e^{\gamma})\, \exp(\gamma(x+1)) \left[1 + \left(\dfrac{\gamma-\log 2}{\log 2}\right)^2\right]^{-1} d\gamma}{\displaystyle\int_{-\infty}^{\infty} \exp(-e^{\gamma})\, \exp(\gamma x) \left[1 + \left(\dfrac{\gamma-\log 2}{\log 2}\right)^2\right]^{-1} d\gamma}$$

Posterior Means under $\pi_1$, $\pi_2$ and $\pi_3$ as follows:

| $\pi$ | 0 | 1 | 2 | 3 | 4 | 5 | 10 | 20 |
|---|---|---|---|---|---|---|---|---|
| $\pi_1$ | .749 | 1.485 | 2.228 | 2.971 | 3.713 | 4.456 | 8.169 | 15.595 |
| $\pi_2$ | .950 | 1.480 | 2.106 | 2.806 | 3.559 | 4.353 | 8.660 | 17.945 |
| $\pi_3$ | .761 | 1.562 | 2.094 | 2.633 | 3.250 | 3.950 | 8.867 | 19.178 |

To see if the choice of prior matters, simply examine the posterior means under the three different priors.

For small or moderate $x$ ($x \leq 10$), there is robustness; the choice of prior does not seem to matter too much. For large values of $x$, the choice does matter. The inference that a conjugate prior obtains then is quite different from what a heavier tailed prior would obtain. It is now clear that there are situations where it does matter what prior one chooses from a class of priors, each of which is considered reasonable given the available prior information.

## Classes of Priors

The goals is to choose a class of priors, $\Gamma$ are,

(i) to ensure that as many 'reasonable' priors as possible are included,

(ii) to try to eliminate 'unreasonable' priors,

(iii) to ensure that $\Gamma$ does not require prior information which is difficult to elicit, and

(iv) to be able to compute measures of robustness without much difficulty.

Example. Suppose $\theta$ is a real-valued parameter, prior beliefs about which indicate that it should have a continuous prior distribution, symmetric about 0 and having the third quartile, $Q_3$, between 1 and 2, Consider, then

$$\Gamma_1 = \{ N(0, \tau^2) : 2.19 < \tau^2 < 8.76 \} \text{ and}$$

$$\Gamma_2 = \{ \text{symmetric priors with } 1 < Q_3 < 2 \}$$

Even though $\Gamma_1$ can be appropriate in some cases, it will mostly be considered "rather small", because it contains only sharp-tailed distribution. On the other hand, $\Gamma_2$ will typically be "too large", containing priors, shapes of some of which will be considered unreasonable.

Starting with $\Gamma_2$ and imposing reasonable

Constraints such as unimodality on the priors can lead to sensible classes such as

$$\Gamma_3 = \{\text{unimodal symmetric priors with } 1 < a_3 < 2\}$$

$$\Gamma_1 \subseteq \Gamma_3 \subseteq \Gamma_2 \ .$$

## Conjugate class

The class considering of conjugate priors is one of the easiest classes of priors to work with. If $X \sim N(\theta, \sigma^2)$ with known $\sigma^2$, the conjugate priors for $\theta$ are the normal priors $N(\mu, \tau^2)$.

So one could consider

$$\Gamma_C = \left\{ N(\mu, \tau^2) : \mu_1 \leq \mu \leq \mu_2 , \tau_1^2 \leq \tau^2 \leq \tau_2^2 \right\}$$

for some specified values of $\mu_1, \mu_2, \tau_1^2,$ and $\tau_2^2$. The advantage with the conjugate class is that posteriors can be calculated in closed form.

In the above case, if $\theta \sim N(\mu, \tau^2)$, then

$$(\theta \mid x = x) \sim N(\mu^*(x), \delta^2) , \text{ where } \delta^2 = \frac{\tau^2 \sigma^2}{\tau^2 + \sigma^2}$$

and $\mu^*(x) = \dfrac{\tau^2}{\tau^2 + \sigma^2} x + \dfrac{\sigma^2}{\tau^2 + \sigma^2} \mu$.

Minimizing and maximizing posterior quantities then becomes any easy task.

## Neighbourhood class

If $\pi_0$ is a single elicited prior, then uncertainity in this elicitation can be modeled using the class

$$\Gamma_N = \{\pi \text{ which are in the neighborhood of } \pi_0\}$$

A natural and well studied class is the $\epsilon$-contamination class,

$$\Gamma_\epsilon = \{\pi : \pi = (1-\epsilon)\pi_0 + \epsilon q, \quad q \in Q\}$$

$\epsilon$ reflecting the uncertainity in $\pi_0$ and $Q$ specifying the contaminations. Some choices for $Q$ are all unimodal distributions with mode $\theta_0$, and all unimodal symmetic distributions with mode $\theta_0$.

The $\epsilon$-contamination class with appropriate choice of $Q$ can provide good robustness.

## Dasity Ratio class

Assuming the existence of densities for all the priors in the class, the density ratio class is defined as $\Gamma_{DR} = \{\pi : L(\theta) \leq \alpha \pi(\theta) \leq U(\theta) \text{ for some } \alpha > 0\}$

$$= \left\{\pi : \frac{L(\theta)}{U(\theta')} \leq \frac{\pi(\theta)}{\pi(\theta')} \leq \frac{U(\theta)}{L(\theta')} \text{ for all } \theta, \theta'\right\}$$

for specified non-negative functions $L$ and $U$.

If we take $L \equiv 1$ and $U \equiv C$ then

$$\Gamma_{DR} = \left\{\pi : C^{-1} \leq \frac{\pi(\theta)}{\pi(\theta')} \leq C, \text{ for all } \theta, \theta'\right\}$$

# Global Measures of Sensitivity

**Example:** Suppose $X_1, X_2, \ldots X_n$ are i.id $N(\theta, \sigma^2)$, with $\sigma^2$ known and $\Gamma$ be all $N(0, \tau^2)$, $\tau^2 > 0$, priors for $\theta$. Then the variation in the posterior mean is

$$\left( \inf_{\tau^2 > 0} E(\theta \mid \bar{x}), \sup_{\tau^2 > 0} E(\theta \mid \bar{x}) \right).$$

Because for fixed $\tau^2$, $E(\theta \mid \bar{x}) = \dfrac{\tau^2}{\tau^2 + \sigma^2} \bar{x}$, this range can easily be seen to be $(0, \bar{x})$ or $(\bar{x}, 0)$ according as $\bar{x} \geq 0$ or $\bar{x} < 0$. If $\bar{x}$ is small in magnitude, this range will be small. Thus the robustness of the procedure of using posterior mean as the Bayes estimate of $\theta$ will depend crucially on the magnitude of the observed value of $\bar{x}$.

As can be seen from the above example, a natural global measure of sensitivity of the Bayesian quantity to the choice of prior is the range of this quantity as the prior varies in the class of priors of interest. Further there are three categories of Bayesian quantities of interest.

(i) Linear functionals of the prior;

$$\rho(\pi) = \int_{\Theta} h(\theta) \pi(d\theta) \qquad \text{where } h \text{ is a given function.}$$

y.

If h is taken to be the likelihood function $l$, we get an important linear functional, the marginal density of of data, i.e.; $m(\pi) = \int_{\Theta} l(\theta) \pi(d\theta)$.

(ii) Ratio of linear functionals of the priors:

$$P(\pi) = \frac{1}{m(\pi)} \int_{\Theta} h(\theta) l(\theta) \pi(d\theta)$$

for some given function h.

(iii) If we take $h(\theta) = \theta$, $P(\pi)$ is the posterior mean.

For $h(\theta) = I_c(\theta)$, the indicator function of the set $C$, we get the posterior probability of $C$.

(iii) Ratio of nonlinear functionals:

$$P(\pi) = \frac{1}{m(\pi)} \int_{\Theta} h(\theta, \phi(\pi)) l(\theta) \pi(d\theta)$$

for some given $h$. For $h(\theta, \phi(\pi)) = (\theta - \mu(\pi))^2$ where $\mu(\pi)$ is the posterior mean, we get

$P(\pi)$ is the posterior variance.

Note that extreme values of linear functionals of the prior as it varies in a class $\Gamma$ are easy to compute if the extreme points of $\Gamma$ can be identified.

## Example

Suppose $X \sim N(\theta, \sigma^2)$, with $\sigma^2$ known and the class $\Gamma$ of interest is

$$\Gamma_{SU} = \{ \text{all symmetric unimodal distributions with mode } \theta_0 \}$$

Then $\phi$ denoting the standard normal density,

$$m(\pi) = \int_{-\infty}^{\infty} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) \pi(\theta) \, d\theta .$$

Note that any unimodal symmetric (about $\theta_0$) density $\pi$ is a mixture of uniform densities symmetric about $\theta_0$. Thus the extreme points of $\Gamma_{SU}$ are $U(\theta_0 - r, \theta_0 + r)$ distribution.

Therefore, $\displaystyle\inf_{\pi \in \Gamma_{SU}} m(\pi) =$

$$= \inf_{r > 0} \frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) d\theta$$

$$= \inf_{r > 0} \frac{1}{2r} \left\{ \Phi\left(\frac{\theta_0 + r - x}{\sigma}\right) - \Phi\left(\frac{\theta_0 - r - x}{\sigma}\right) \right\} .$$

$$\sup_{\pi \in \Gamma_{SU}} m(\pi) = \sup_{r > 0} \frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) d\theta$$

$$= \sup_{r > 0} \frac{1}{2r} \left\{ \Phi\left(\frac{\theta_0 + r - x}{\sigma}\right) - \Phi\left(\frac{\theta_0 - r - x}{\sigma}\right) \right\}$$

## Lemma

Suppose $C_T$ is a set of probability measures on the real line, given by $C_T = \{ v_t : t \in T \}$, $T \subset \mathbb{R}^d$, and let $\mathcal{C}$ be the convex hull of $C_T$. Further suppose $h_1$ and $h_2$ are real-valued functions defined on $\mathbb{R}$ such that $\int |h_1(x)| \, dF(x) < \infty$ for all $F \in \mathcal{C}$, and $K + h_2(x) > 0$ for all $x$, for some constant $K$. Then, for any $k$,

$$\sup_{F \in \mathcal{C}} \frac{k + \int h_1(x) \, dF(x)}{K + \int h_2(x) \, dF(x)} = \sup_{t \in T} \frac{k + \int h_1(x) \, v_t(dx)}{K + \int h_2(x) \, v_t(dx)}$$

$$\inf_{F \in \mathcal{C}} \frac{k + \int h_1(x) \, dF(x)}{K + \int h_2(x) \, dF(x)} = \inf_{t \in T} \frac{k + \int h_1(x) \, v_t(dx)}{K + \int h_2(x) \, v_t(dx)}$$

**Proof** Because $\int h_1(x) \, dF(x) = \int h_1(x) \int_T v_t(dx) \, \mu(dt)$, for some probability measure $\mu$ on $T$, using Fubini's theorem,

$$k + \int h_1(x) \, dF(x) = \int (k + h_1(x)) \int_T v_t(dx) \, \mu(dt)$$

$$= \int_T \left( \int (k + h_1(x)) \, v_t(dx) \right) \mu(dt)$$

$$= \int_T \left[ \left( \frac{\int (k + h_1(x)) \, v_t(dx)}{\int (K + h_2(x)) \, v_t(dx)} \right) \int (K + h_2(x)) \, v_t(dx) \right] \mu(dt)$$

$$\leq \left( \sup_{t \in T} \frac{\int (k + h_1(x)) \, \nu_t(dx)}{\int (k + h_2(x)) \, \nu_t(dx)} \right) \left( K + \int h_2(x) \, dF(x) \right) \; .$$

Therefore,

$$\sup_{F \in \mathcal{C}} \frac{k + \int h_1(x) \, dF(x)}{K + \int h_2(x) \, dF(x)} \leq \sup_{t \in T} \frac{\int (k + h_1(x)) \, \nu_t(dx)}{\int (K + h_2(x)) \, \nu_t(dx)}$$

However, because $\mathcal{C} \supset C_T$,

$$\sup_{F \in \mathcal{C}} \frac{k + \int h_1(x) \, dF(x)}{K + \int h_2(x) \, dF(x)} \geq \sup_{t \in T} \frac{\int (k + h_1(x)) \, \nu_t(dx)}{\int (K + h_2(x)) \, \nu_t(dx)}$$

Hence the proof for the supremum, and the proof for the infimum is similar.

**Therem** Consider the class $\Gamma_{SU}$ of all symmetric unimodal prior distributions with mode $\theta_0$. Then it follows that

$$\sup_{\pi \in \Gamma_{SU}} E^{\pi}(g(\theta) | x) = \sup_{r > 0} \frac{\frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} g(\theta) \, f(x | \theta) \, d\theta}{\frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} f(x | \theta) \, d\theta},$$

$$\inf_{\pi \in \Gamma_{SU}} E^{\pi}(g(\theta) | x) = \inf_{r > 0} \frac{\frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} g(\theta) \, f(x | \theta) \, d\theta}{\frac{1}{2r} \int_{\theta_0 - r}^{\theta_0 + r} f(x | \theta) \, d\theta}$$

**Proof** Note that $E^{\pi}(g(\theta) | x) = \dfrac{\int g(\theta) f(x | \theta) \, d\pi(\theta)}{\int f(x | \theta) \, d\pi(\theta)}$,

where $f(x | \theta)$ is the density of the data $x$.

Now using the above lemma and recalling that any unimodal symmetric distribution is a mixture of symmetric uniform distribution, we get the result.

**Example** Suppose $X \mid \theta \sim N(\theta, \sigma^2)$ and robustness
· of the posterior mean w.r.t. $\Gamma_{SU}$ is of interest.
Then, range of posterior mean over this class can be
easily computed using the previous theorem. Thus

$$\sup_{\pi \in \Gamma_{SU}} E^{\pi}(\theta \mid x) = \sup_{\gamma > 0} \frac{\frac{1}{2\gamma} \int_{\theta_0 - \gamma}^{\theta_0 + \gamma} \frac{\theta}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) d\theta}{\frac{1}{2\gamma} \int_{\theta_0 - \gamma}^{\theta_0 + \gamma} \frac{1}{\sigma} \phi\left(\frac{x - \theta}{\sigma}\right) d\theta}$$

$$= x + \sup_{\gamma > 0} \frac{\phi\left(\frac{\theta_0 - \gamma - x}{\sigma}\right) - \phi\left(\frac{\theta_0 + \gamma - x}{\sigma}\right)}{\Phi\left(\frac{\theta_0 + \gamma - x}{\sigma}\right) - \Phi\left(\frac{\theta_0 - \gamma - x}{\sigma}\right)}$$

$$\inf_{\pi \in \Gamma_{SU}} E^{\pi}(\theta \mid x) = x + \inf_{\gamma > 0} \frac{\phi\left(\frac{\theta_0 - \gamma - x}{\sigma}\right) - \phi\left(\frac{\theta_0 + \gamma - x}{\sigma}\right)}{\Phi\left(\frac{\theta_0 + \gamma - x}{\sigma}\right) - \Phi\left(\frac{\theta_0 - \gamma - \gamma}{\sigma}\right)}$$

**Example** Suppose $X \mid \theta \sim N(\theta, \sigma^2)$ and it is of
interest to test $H_0: \theta \leq \theta_0$ vs $H_1: \theta > \theta_0$. Again,
suppose that $\Gamma_{SU}$ is the class of priors to be
considered and robustness of this class is to be
examined. Because

$$P^{\pi}(H_0 \mid x) = P^{\pi}(\theta \leq \theta_0 \mid x)$$

$$= \frac{\int_{-\infty}^{\infty} I_{(-\infty, \theta_0]}(\theta) f(x \mid \theta) \, d\pi(\theta)}{\int_{-\infty}^{\infty} f(x \mid \theta) \, d\pi(\theta)}$$

appling the previous theorem, we get

$$\sup_{\pi \in \Gamma_{su}} P^\pi(H_0|x) = \sup_{r>0} \frac{\frac{1}{2r} \int_{\theta_0-r}^{\theta_0} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) d\theta}{\frac{1}{2r} \int_{\theta_0-r}^{\theta_0+r} \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) d\theta}$$

$$= \sup_{r>0} \frac{\Phi\left(\frac{\theta_0-x}{\sigma}\right) - \Phi\left(\frac{\theta_0-r-x}{\sigma}\right)}{\Phi\left(\frac{\theta_0+r-x}{\sigma}\right) - \Phi\left(\frac{\theta_0-r-x}{\sigma}\right)}$$

and similarly,

$$\inf_{\pi \in \Gamma_{su}} P^\pi(H_0|x) = \inf_{r>0} \frac{\Phi\left(\frac{\theta_0-x}{\sigma}\right) - \Phi\left(\frac{\theta_0-r-x}{\sigma}\right)}{\Phi\left(\frac{\theta_0+r-x}{\sigma}\right) - \Phi\left(\frac{\theta_0-r-x}{\sigma}\right)}$$

It can be seen that the above bounds are, respectively, $0.5$ and $d$, where $d = \Phi\left(\frac{x-\theta_0}{\sigma}\right)$, the $P$-value.

We shall now consider the density-ratio class that was metioned earlier is given by

$$\Gamma_{DR} = \left\{ \pi : L(\theta) \le \alpha \pi(\theta) \le U(\theta) \text{ for some } \alpha > 0 \right\}$$

for specified non-negative functions $L$ and $U$.

For $\pi \in \Gamma_{DR}$ and any real valued $\pi$-integrable function $h$ on the parameter space $\Theta$, let

$$\pi(h) = \int_\Theta h(\theta) \pi(d\theta). \quad \text{and} \quad \text{let } h \equiv h^+ - h^-,$$

## Theorem

For $U$-integrable functions $h_1$ and $h_2$, with $h_2$ positive a.s. w.r.t. all $\pi \in \Gamma_{DR}$,

$$\inf_{\pi \in \Gamma_{DR}} \frac{\pi(h_1)}{\pi(h_2)} \quad \text{is the unique solution } \lambda \text{ of}$$

$$U(h_1 - \lambda h_2)^- + L(h_1 - \lambda h_2)^+ = 0$$

$$\sup_{\pi \in \Gamma_{DR}} \frac{\pi(h_1)}{\pi(h_2)} \quad \text{is the unique solution } \lambda \text{ of}$$

$$U(h_1 - \lambda h_2)^+ + L(h_1 - \lambda h_2)^- = 0$$

**Proof** let $\lambda_0 = \inf_{\pi \in \Gamma_{DR}} \frac{\pi(h_1)}{\pi(h_2)}$, $c_1 = \inf_{\pi \in \Gamma_{DR}} \pi(h_2)$

and $c_2 = \sup_{\pi \in \Gamma_{DR}} \pi(h_2)$. Then $0 < c_1 < c_2 < \infty$,

and $|\lambda_0| < \infty$. Because $U(h_1 - \lambda h_2)^- +$

$L(h_1 - \lambda h_2)^+ = \inf_{\pi \in \Gamma_{DR}} \pi(h_1 - \lambda h_2)$ for any $\lambda$,

note that $\lambda_0 \geq \lambda$ iff $U(h_1 - \lambda h_2)^- + L(h_1 - \lambda h_2)^+ \geq 0$

However, $\lambda_0 > \lambda$ iff $U(h_1 - \lambda h_2)^- + L(h_1 - \lambda h_2)^+ > 0$.

A si Hence $\lambda_0 = \lambda$ iff $U(h_1 - \lambda h_2)^- + L(h_1 - \lambda h_2)^+ = 0$

A similar argument for the supremum.

## Example

Suppose $X \sim N(\theta, \sigma^2)$, with $\sigma^2$ known. Consider the class $\Gamma_{DR}$ with $L$ being the Lebesgue measure and $U = kL$, $k > 1$. Because the posterior mean is

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

we have that $\displaystyle\inf_{\pi \in \Gamma_{DR}} E^{\pi}(\theta|x)$ is the unique

solution $\lambda$ of

$$k \int_{-\infty}^{\lambda} (\theta - \lambda) f(x|\theta) \, d\theta + \int_{\lambda}^{\infty} (\theta - \lambda) f(x|\theta) \, d\theta = 0$$

and similarly, $\displaystyle\sup_{\pi \in \Gamma_{DR}} E^{\pi}(\theta|x)$ is the unique

solution $\lambda$ of $\displaystyle\int_{-\infty}^{\lambda} (\theta - \lambda) f(x|\theta) \, d\theta + k \int_{\lambda}^{\infty} (\theta - \lambda) f(x|\theta) \, d\theta = 0$

Noting that $f(x|\theta) = \frac{1}{\sigma} \phi\left(\frac{x-\theta}{\sigma}\right) = \frac{1}{\sigma} \phi\left(\frac{\theta-x}{\sigma}\right)$,

and letting $\lambda_1$ be the minimum and $\lambda_2$ the maximum, the above equations may be rewritten as

$$(k-1)\left[\left(\frac{\lambda_1-x}{\sigma}\right)\Phi\left(\frac{\lambda_1-x}{\sigma}\right)+\phi\left(\frac{\lambda_1-x}{\sigma}\right)\right]=\frac{\lambda_1-x}{\sigma},$$

$$(k-1)\left[\left(\frac{\lambda_2-x}{\sigma}\right)\Phi\left(\frac{\lambda_2-x}{\sigma}\right)+\phi\left(\frac{\lambda_2-x}{\sigma}\right)\right]=k\left(\frac{\lambda_2-x}{\sigma}\right).$$

Now let $k\left(\frac{\lambda_1-x}{\sigma}\right)=\gamma$, Then $\lambda_2=x+\sigma\frac{\gamma}{k}$

Put $\lambda_0=x-\sigma\frac{\gamma}{k}$ or $\frac{\lambda_0-x}{\sigma}=-\frac{\gamma}{k}$,

Then we see from the second equation above that

$$(k-1)\left[\left(\frac{\lambda_0-x}{\sigma}\right)\Phi\left(\frac{\lambda_0-x}{\sigma}\right)+\phi\left(\frac{\lambda_0-x}{\sigma}\right)\right]$$

$$=(k-1)\left[-\frac{\gamma}{k}\Phi\left(-\frac{\gamma}{k}\right)+\phi\left(-\frac{\gamma}{k}\right)\right]$$

$$=(k-1)\left[-\frac{\gamma}{k}\left(1-\Phi\left(\frac{\gamma}{k}\right)\right)+\phi\left(\frac{\gamma}{k}\right)\right]$$

$$=(k-1)\left[\frac{\gamma}{k}\Phi\left(\frac{\gamma}{k}\right)+\phi\left(\frac{\gamma}{k}\right)\right]-(k-1)\frac{\gamma}{k}$$

$$=0$$

implying that once $\lambda_2$ is obtained, say

$\lambda_2=x+\sigma\cdot\frac{\gamma}{k}$, the solution for $\lambda_1$ is

simply $x-\sigma\frac{\gamma}{k}$.

**Different Methods of Construction of Objective Priors:**

To construct objective priors under general regularity conditions, one may do one of the following things.

1. Define a uniform distribution that takes into account the geometry of the parameter space.

2. Minimize a suitable measure of information in the prior.

3. Choose a prior with some form of frequentist ideas because a prior with little information should lead ~~lead~~ to inference that is similar to frequentist inference.

Although the usual uniform prior $\pi(\theta) = c$ has come in for a lot of criticism, these criticisms help one understand the motivation behind (1) and (2). The fact is that both (1) and (2) lead to the Jeffreys prior,

namely, $\pi(\theta) = [\det(I_{ij}(\theta))]^{1/2}$

where $(I_{ij}(\theta))$ is the Fisher information matrix. In the one-dimensional case (3) also leads to the Jeffreys prior.

# Uniform Distribution and Its Criticisms

The first objective prior ever to be used is the uniform distribution over a bounded interval.

A common argument, based on "ignorance", seems to have been that if we know nothing about $\theta$, why should we attach more density to one point than another? ~~The argumen~~

A second argument is that the uniform maximizes the shannon entropy.

The principle of ignorance has been criticized by many statisticians. 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 for $\eta$ also. So the principle of ignorance applied to $\eta$ will imply our prior for $\eta$ is uniform on $\psi(\Theta)$ just as it had led to a uniform prior for $\theta$. But this leads to a contradiction. To see this suppose $\psi$ is a differentiable and $p(\eta) = c$ on $\psi(\Theta)$.

Then the prior $p^*(\theta)$ for $\theta$ is

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

Which is not a constant in general.

This argument also leads to an invariance principle. Suppose we have an algorithm that produces noninformative 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 noninformative prior should be invariant under one-to-one differential transformation.

The second argument in favour of the uniform, based on shannon entropy, is also flawed. Shannon derived a measure of entropy in the finite discrete case from certain natural axioms. His entropy is $H(p) = -\sum_{i=1}^{m} p_i \log p_i$ which is maximized by the discrete uniform, i.e., at $p = \left(\frac{1}{m}, \cdots, \frac{1}{m}\right)$.

Entropy is a measure of the amount of uncertainty about the outcomes of the experiment.

A prior that maximizes this will maximize the uncertainity, so it is a noninformative prior. Because such a prior should minimize the information, we take negative of entropy as information.

Shannon's entropy is a natural measure in the discrete case and the discrete uniform appears to be the right noninformative prior. The continuous case is an entirely different matter. Shannon pointed out that for a density $p$,

$$H(p) = -\int (\log p(x)) p(x) \, dx$$

is unsatisfactory, clearly it is not derived from axioms, it is not invariant under one-to-one transformation, also, it depends on the measure $\mu(x) dx$ w.r.t which the density $p(x)$ is taken. Note also that the measure $\mu(x) dx$ is not non-negative. Just take $\mu(x) = 1$ and take $p(x)$ = uniform on $[0, c]$, then $H(p) > 0$ iff $c > 1$.

Finally, if the density is take w.r.t $\mu(x) dx$, then it is easy to verify that the density is

$$\frac{p}{\mu} \quad \text{and} \quad H(p) = -\int \left( \log \frac{p(x)}{\mu(x)} \right) \frac{p(x)}{\mu(x)} \mu(x) \, dx$$

is maximized at $p = \mu$, i.e., the entropy is maximum at the arbitrary $\mu$. arbitrary $\mu$. For all these reasons, we do not think $H(p)$ Is the right entropy to maximize. However, $H(p)$ serves a useful purpose when we have partial information.

# Jeffreys Priors as a Uniform Distribution

Suppose $\Theta = \mathbb{R}^d$ and $I(\theta) = (I_{ij}(\theta))$ is the $d \times d$ Fisher information matrix. We assum $I(\theta)$ is positive definite for all $\theta$. Rao had proposed the Riemannian metric $\rho$ related to $I(\theta)$ by

$$\rho(\theta, \theta + d\theta) = \sum_{i,j} I_{i,j}(\theta)\, d\theta_i\, d\theta_j\, (1 + o(1))$$

It is known that this is the unique Riemannian metric that transforms suitably under one-one differentiable transformations on $\Theta$. Notice that in general $\Theta$ does not inherit the usual Euclidean metric that goes with the improper uniform distribution over $\mathbb{R}^d$.

Fix a $\theta_0$ and let $\psi(\theta)$ be a smooth one-to-one transformation such that the information matrix $I^\psi = \left[ E_\psi \left( \dfrac{\partial \log p}{\partial \psi_i} \cdot \dfrac{\partial \log p}{\partial \psi_j} \right) \right]$ is the identity matrix $I$ at $\psi_0 = \psi(\theta_0)$. This implies the local geometry in the $\psi$-space around $\psi_0$ is Euclidean and hence $d\psi$ is a suitable uniform distribution there. If we now lift this back to the $\theta$-space by using the Jacobian of transformation and the simple fact

$$\left(\left|\frac{\partial \theta_j}{\partial \psi_i}\right|\right)\left(I_{ij}(\theta)\right)\left(\left|\frac{\partial \theta_j}{\partial \psi_i}\right|\right)' = I^{\psi_0} = I,$$

we get the Jeffreys prior in the $\theta$-space,

$$d\psi = \left\{ \det\left|\frac{\partial \theta_i}{\partial \psi_j}\right|\right\}^{-1} d\theta = \left\{\det\left[I_{ij}(\theta)\right]\right\}^{1/2} d\theta.$$

A similar method present an alternative construction where one takes a compact subset of the parameter space and approximates this by a finite set of points in the so-called Hellinger metric

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

where $P_\theta$ and $P_{\theta'}$ are the densities of $P_\theta$ and $P_{\theta'}$. One then puts a discrete uniform distribution on the approximating finite set of points, and lets the degree of approximation tends to zero. Then the corresponding discrete uniforms converge weakly to the Jeffreys distribution.

# Jeffreys Prior as a Minimizer of Information

Let the Shannon entropy associated with a random variable or vector $Z$ be denoted by

$$H(Z) = H(p) = - E_p(\log p(Z))$$

where $p$ is the density of $Z$. Let $X = (x_1, x_2, \cdots x_n)$ have density or probability function $p(x|\theta)$ where $\theta$ has prior distribution density $p(\theta)$. We assume $X_1, X_2, \cdots X_n$ are i.i.d and conditions for asymptotic nomality of posterior $p(\theta|x)$ hold. We know that $H(p)$ is not a good measure of entropy and $-H(p)$ not a good measure of information if $p$ is a density. Bernardo suggested. suggested a Kullback-Leibler divergence between prior and posterior, namely,

$$J(p(\theta), X) = E\left\{ \log \frac{p(\theta|X)}{p(\theta)} \right\}$$

$$= \int_{\Theta} \left\{ \int_X \left[ \int_{\Theta} \log \left\{ \frac{p(\theta'|x)}{p(\theta')} \right\} p(\theta'|x) d\theta' \right] p(x|\theta) dx \right\} p(\theta) d\theta$$

is a better measure of entropy and $-J$ a better measure of information in the prior.

To get a feeling for this, notice that if the prior is nearly degenerate, at pay some $\theta_0$, so will be the posterior. This would imply $J$ is nearly zero. On the other hand if $P(\theta)$ is rather diffuse, $P(\theta|x)$ will be differ a lot form $P(\theta)$, at least for moderate or large $n$, because $p(\theta|x)$ would be quite peaked. In fact $P(\theta|x)$ would be approximately normal with mean $\hat{\theta}$ and variance of the order $O(\frac{1}{n})$. The substantial difference between prior and posterior would be reflected by a large value of $J$.

To sum up $J$ is small when $p$ is nearly degenerate and large when $p$ is diffuse, i.e, $J$ captures how diffuse is the prior. It therefore makes sense to minimize $J$ with respect to the prior.

After a lot of calculation one can show that the Jeffreys prior minimize $J$.

## Measure of entropy or information

Shannon introduced missing information in the context of a noisy channel. Any channel has a source that produces (say, per second) messages X with p.m.f $P_X(x)$ and entropy

$$H(X) = - \sum_x p_x(x) \log P_X(x) .$$

A channel will have an output Y (per second) with entropy $H(Y) = - \sum_y P_Y(y) \log P_Y(y)$.

If the channel is noiseless, then $H(Y) = H(X)$.

If the channel is noisy, Y given X is still random. Let $P(x,y)$ denote their joint p.m.f.

The joint entropy is

$$H(x,y) = - \sum_{x,y} p(x,y) \log P(x,y).$$

Clearly, $H(x,y) = H(x) + H_x(Y) = H(Y) + H_Y(x)$

~~and $H(x,y)$~~ $H_Y(x)$ is called the equivocation or average ambiguity about input X given only output Y. It is the information about input X that is received given the output Y.

It is the amount of additional information that must be supplied per unit time at the receiving end to correct the received message.

Thus $H_Y(x)$ is the missing information. So amount of information produced in the channel (per unit time) is $H(x) - H_Y(x)$

which may be shown to be non-negative by Shannon's basic results,

$$H(x) + H(Y) \geq H(x, Y) = H(Y) + H_Y(x).$$

In statistical problems, we take $X$ to be $\theta$ and $Y$ to be observation vector $x$. Then $H(\theta) - H_x(\theta)$ is the same measure as before, namely

$$E \left( \log \frac{p(\theta/x)}{p(\theta)} \right)$$

The maximum of $H(x) - H_Y(x)$ with respect to the source, i.e., w.r.t. $p(x)$ is what Shannon calls the capacity of the channel. Over compact rectangles, the Jeffreys prior is this maximizing distribution for the statistical channel.

Hypothesis testing and Model Selection

Suppose X having density $f(x|\theta)$ is observed, with $\theta$ being an unknown element of the parameter space $\Theta$. Suppose that we are interested in comparing two models $M_0$ and $M_1$, which is given by

$M_0$: X has density $f(x|\theta)$ where $\theta \in \Theta_0$.

$M_1$: X has density $f(x|\theta)$ where $\theta \in \Theta_1$.

For $i = 0, 1$ let $g_i(\theta)$ be the prior density of $\theta$, conditional on $M_i$ being the true model. Then, to compare models $M_0$ and $M_1$ on the basis of a random sample $x = (x_1, x_2, \cdots x_n)$ one would use the Bayes factor $B_{01}(x) = \dfrac{m_0(x)}{m_1(x)}$

where $m_i(x) = \displaystyle\int_{\Theta_i} f(x|\theta)\, g_i(\theta)\, d\theta$, $i = 0, 1$

We also use the notation $BF_{01}$ for the Bayes factor, $B_0$ which is the ratio of posterior odds ratio of the hypotheses to the corresponding prior odds ratio. Therefore, if the prior probabilities of the hypotheses, $\pi_0 = P^{\pi}(M_0) = P^{\pi}(\Theta_0)$ and $\pi_1 = P^{\pi}(M_1) = P^{\pi}(\Theta_1) = 1 - \pi_0$ are specified,

then
$$P(M_0 | x) = \left\{ 1 + \frac{1 - \pi_0}{\pi_0} B_{01}^{-1}(x) \right\}^{-1}.$$

Thus, if conditional prior densities $g_0$ and $g_1$ can be specified, one should simply use the Bayes factor $B_{01}$ for model selection.

If, further $\pi_0$ is also specified, the posterior odds ratio of $M_0$ to $M_1$ can also be utilized. However, these computations may not always be easy to perform, even when the required prior ingredients are fully specified. A possible solution is the use of BIC as an approximation to a Bayes factor.

Example Consider the problem that is usually called non parametric regression. Independent responses $Y_i$ are observed along with covariates $x_i$, $i = 1, 2, \ldots n$.

The model of interest is

$$Y_i = g(x_i) + \epsilon_i, \quad i = 1, 2, \ldots n$$

where $\epsilon_i$ are iid $N(0, \sigma^2)$ errors with unknown error variance $\sigma^2$. The function $g$ is called the regression function.

In linear regression, $g$ is a priori assumed to be ~~fully unknown also~~. linear in a set of finite regression coefficients. In general, $g$ can be assumed to be fully unknown also. Now, if model selection involves choosing $g$ from two different fully nonparametric classes of regression functions, this becomes a very difficult problem. various simplifications including reducing $g$ to be semi-parametric have been studied.

Consider a different model checking problem, that of testing for normality. In its simplest ~~for~~ form, the problem can be stated as checking whether a given random sample $X_1, X_2, \cdots, X_n$ arose from a population having the normal distribution. We may write it as

$M_0$: $X$ is $N(\mu, \sigma^2)$ with arbitrary $\mu$ and $\sigma^2 > 0$.

$M_1$: $X$ does not have the normal distribution.

However, this looks quite different, because $M_1$ does not constitute a parametric alternative. Hence it is not clear how to use Bayes factors or posterior odds ratios here for model checking.

# P-value and Posterior Probability of $H_0$ as measures of Evidence against the Null.

One particular tool from classical statistics that is very widely used in applied sciences for model checking or hypothesis testing is the P-value. It is also happens to be one of the concepts that is highly misunderstood and misused. The basic idea behind the definition of P-value is,

It is the probability of under a null hypothesis of obtaining a value of a test statistic that is at least as extreme as that observed in the sample data.

Suppose that it is desired to test

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

and that a classical significance test is available and is based on a test statistic $T(x)$, large values of which are deemed to provide evidence against the null hypothesis. If data $X = x$ is observed, with corresponding $t = T(x)$, the P-value then is

$$\alpha = P_{\theta_0}(T(x) \geqslant T(x))$$

## Example

Suppose we observed $X_1, X_2, \ldots, X_n$ i.i.d from $N(\theta, \sigma^2)$; where $\sigma^2$ is known. Then $\bar{X}$ is sufficient for $\theta$ and it has the $N(\theta, \sigma^2/n)$ distribution.

Noting that $T = T(\bar{x}) = |\sqrt{n}(\bar{x} - \theta_0)/\sigma|$ is a natural test statistics to test

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

one obtains the usual P-value as $\alpha = 2[1 - \Phi(t)]$ where $t = |\sqrt{n}(\bar{x} - \theta_0)/\sigma|$ and $\Phi$ is the standard normal cumulative distribution function.

Fisher meant P-value to be used informally as a measure of degree of surprise in the data relative to $H_0$. This use of P-value as a post-experimental or conditional measure of statistical evidence seems to have some intuitive justification.

## Bounds on Bayes Factors and Posterior Probabilities

Consider the following example where P-values and the posterior probabilities are very different.

We observe $\bar{x} \sim N(\theta, \sigma^2/n)$, with known $\sigma^2$. Upon using $T = |\sqrt{n}(\bar{x} - \theta_0)/\sigma|$ as the test statistics to test

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

Then the P-values comes out to be $\alpha = 2[1 - \Phi(t)]$

On the set $\{\theta \neq \theta_0\}$, let $\theta$ have the density $g_1$ of $N(\mu, \tau^2)$. Then we have

$$B_{01} = \sqrt{1 + \rho^{-2}} \; \exp\left\{-\frac{1}{2}\left[\frac{(t - \rho\eta)^2}{(1 + \rho^2)} - \eta^2\right]\right\}$$

where $\rho = \sigma/(\sqrt{n}\,\tau)$ and $\eta = (\theta_0 - \mu)/\tau$.

Now, if we choose $\mu = \theta_0$, $\tau = \sigma$ and $\pi_0 = \frac{1}{2}$

we get $$B_{01} = \sqrt{1 + \rho^{-2}} \; \exp\left\{-\frac{1}{2}\left[\frac{t^2}{1 + \rho^2}\right]\right\}$$

For various values of $t$ and $n$, the different measures of evidence, $\alpha = $ P-value, $B = $ Bayes factor, and $P = P(H_0|x)$ are displayed is as follows.

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | $n$ | | | | |
| | | 1 | | 5 | | 10 | | 20 | | 50 | | 100 |
| | | | | | | | | | | | | |
| $t$ | $\alpha$ | B | P | B | P | B | P | B | P | B | P | B | P |
| 1.645 | 0.10 | .72 | .42 | | | | | | | | | | |
| 1.960 | 0.05 | .54 | .35 | | | | | | | | | | |
| 2.576 | 0.01 | .27 | .21 | | | | | | | | | | |
| 3.291 | .001 | .10 | .09 | | | | | | | | | | |

It may be noted that the posterior probability of $H_0$ varies between 4 to 50 times the corresponding P-value which is an indication of how different these two measures of evidence can be.

# Bayesian Computation

Bayesian analysis requires lots of computation of expectations and quantiles of probability distributions that arise as posterior distributions. Modes of the densities of such distributions are also sometimes used. The standard Bayes estimate is the posterior mean, which is also the Bayes rule under the squared error loss. Its accuracy is assessed using the posterior variance, which is again an expected value. Posterior median is sometimes utilized, and to provide Bayesian credible regions, quantilies of posterior distributions are needed. If conjugate priors are not used, as is mostly the case these days, posterior distributions will not be standard distributions and hence the required Bayesian quantities (i.e., posterior quantities of inferential interest) cannot be computed in closed from. The special techniques are needed for Bayesian computations.

**Example 1** Suppose $X$ is $N(\theta, \sigma^2)$ with known $\sigma^2$ and Cauchy $(\mu, \tau)$ prior on $\theta$ is considered appropriate from a robustness considerations.

Then $\pi(\theta|x) \propto \exp\left(-(\theta-x)^2/(2\sigma^2)\right)\left(\tau^2 + (\theta-\mu)^2\right)^{-1}$

and hence the posterior mean and variance are

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

and $V^{\pi}(\theta|x) = \dfrac{\int_{-\infty}^{\infty} \theta^2 \exp\left(-\frac{(\theta-x)^2}{2\sigma^2}\right)\left(\tau^2 + (\theta-\mu)^2\right)^{-1} d\theta}{\int_{-\infty}^{\infty} \exp\left(-\frac{(\theta-x)^2}{2\sigma^2}\right)\left(\tau^2 + (\theta-\mu)^2\right)^{-1} d\theta} - \left(E^{\pi}(\theta|x)\right)^2$

Note that the above integrals cannot be computed in closed form, but various numerical integration techniques such as IMSL routines or Gaussian quadrature can be efficiently used to obtain very good approximations of these.

Example 2.    Suppose $X_1, X_2, \ldots, X_K$ are independent Poisson counts with $X_i \sim \text{Poisson}(\theta_i)$. $\theta_i$ are a priori considered related, and a joint multivariate normal prior distribution on their logarithm is assumed. Specifically, let $v_i = \log(\theta_i)$ be the $i$th element of $\underline{v}$ and suppose

$$\underline{v} \sim N_K\left(\mu\underline{1}, \ \tau^2\{(1-\rho)I_K + \rho\underline{1}\underline{1}'\}\right)$$

where $\underline{1}$ is the $k$-vector with all elements being 1, and $\mu$, $\tau^2$ and $\rho$ are known constants.

Then, because

$$f(\underline{x}|\underline{v}) = \exp\left(-\sum_{i=1}^{K}\{e^{v_i} - v_i x_i\}\right)\Big/ \prod_{i=1}^{K} x_i!$$

and $\pi(\underline{v}) \propto \exp\left(-\frac{1}{2\tau^2}(\underline{v} - \mu\underline{1})'\left((1-\rho)I_K + \rho\underline{1}\underline{1}'\right)^{-1}(\underline{v}-\mu\underline{1})\right)$

we have that $\pi(\underline{v}|\underline{x}) \propto$

$$\exp\left\{-\sum_{i=1}^{K}\{e^{v_i} - v_i x_i\} - \frac{1}{2\tau^2}(\underline{v}-\mu\underline{1})'\left((1-\rho)I_K + \rho\underline{1}\underline{1}'\right)^{-1}(\underline{v}-\mu\underline{1})\right\}$$

There fore, if the posterior mean of $\theta_j$ is of interest, we need to compute :

$$E^{\pi}(\theta_j \mid x) = E^{\pi}(exp(v_j) \mid x) = \frac{\int_{R^k} exp(v_j)\, g(\underline{v} \mid \underline{z})\, d\underline{v}}{\int_{R^k} g(\underline{v} \mid \underline{z})\, d\underline{v}}$$

where $g(\underline{v} \mid \underline{z}) = $

$$exp\left\{ - \sum_{i=1}^{k} \{ e^{v_i} - v_i x_i \} - \frac{1}{2\tau^2} (\underline{v} - \mu\underline{1})' ((1-\rho) I_k + \rho \underline{1}\underline{1}')^{-1} (\underline{v} - \mu\underline{1}) \right\}$$

This is a ratio of two k-dimensional integrals and as k grows; the integrals become less and less easy to work with. Numerical integration techniques fail to be an efficient technique in this case. In fact, numerical integration techniques are presently not preferred except for single and two-dimensional integrals.

## Monte Carlo Sampling

Consider an expectation that is not available in closed ~~from~~ form. An alternative to numerical integration or analytic approximation to compute this is statistical sampling. 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 the sample proportion. The law of large numbers guarantees that the estimates so obtained will be good provided the sample is large enough.

specifically, let $f$ be a probability density function (or a mass function) and suppose the quantity of interest is a finite expectation of the form

$$E_f \, h(\underline{x}) = \int_X h(\underline{x}) \, f(\underline{x}) \, d\underline{x}$$

(or the corresponding sum in the discrete case).

If i.i.d. observations $x_1, x_2, \ldots$ can be generated from the density $f$, then

$$\bar{h}_m = \frac{1}{m} \sum_{i=1}^{m} h(\underline{x}_i)$$

converges in probability (or even almost surely) to $E_f \, h(\underline{x})$. This justifies using $\bar{h}_m$ as an approximation for $E_f h(\underline{x})$ for large $m$. To provide a measure of accuracy or the extent of error in the approximation, we can again use a statistical technique and compute the standar error. If $Var_f \, h(\underline{x})$ is finite, then $Var_f(\bar{h}_m) = Var_f \, h(\underline{x})/m$.

Further $Var_f \, h(\underline{x}) = E_f \, h^2(\underline{x}) - (E_f \, h(\underline{x}))^2$ can be estimated by

$$S_m^2 = \frac{1}{m} \sum_{i=1}^{m} (h(\underline{x}_i) - \bar{h}_m)^2$$

and hence the standard error of $\bar{h}_m$ can be estimated by

$$\frac{1}{\sqrt{m}} S_m = \frac{1}{m} \left( \sum_{i=1}^{m} (h(\underline{x}_i) - \bar{h}_m)^2 \right)^{1/2}$$

Confidence intervals for $E_f h(\underline{x})$ can be also be provided using the CLT. Because

$$\frac{\sqrt{m} \left( \bar{h}_m - E_f h(\underline{x}) \right)}{S_m} \xrightarrow{d} N(0,1) \qquad \text{as } m \to \infty$$

in distribution, $\left( \bar{h}_m - z_{\alpha/2} S_m/\sqrt{m}, \; \bar{h}_m + z_{\alpha/2} S_m/\sqrt{m} \right)$

can be used as an approximate $100(1-\alpha)\%$ confidence interval for $E_f h(\underline{x})$, with $z_{\alpha/2}$ denoting the $100(1-\alpha/2)\%$ quantile of standard normal.

The above discussion suggests that if we want to approximate the posterior mean, we could try to generate i.i.d. observations from the posterior distribution and consider the mean of this sample. This is rarely useful because most often the posterior distribution will be a non-standard distribution which may not easily allow sampling from it.

**Example** Recall the example 1.

$$E^{\pi}(\theta \mid x) = \frac{\int_{-\infty}^{\infty} \theta \exp\left(-\frac{(\theta-x)^2}{2\sigma^2}\right)\left(\tau^2 + (\theta-\mu)^2\right)^{a-1} d\theta}{\int_{-\infty}^{\infty} \exp\left(-\frac{(\theta-x)^2}{2\sigma^2}\right)\left(\tau^2 + (\theta-\mu)^2\right)^{-1} d\theta}$$

$$= \frac{\int_{-\infty}^{\infty} \theta \left\{\frac{1}{\sigma} \phi\left(\frac{\theta-x}{\sigma}\right)\right\}\left(\tau^2 + (\theta-\mu)^2\right)^{-1} d\theta}{\int_{-\infty}^{\infty} \left\{\frac{1}{\sigma} \phi\left(\frac{\theta-x}{\sigma}\right)\right\}\left(\tau^2 + (\theta-\mu)^2\right)^{-1} d\theta}$$

where $\phi$ denotes the density of standard normal.
Thus $E^{\pi}(\theta \mid x)$ is the ratio of expectation of $h\theta$
$h(\theta) = \theta / (\tau^2 + (\theta-\mu)^2)$ to that of $h(\theta) = 1/(\tau^2 + (\theta-\mu)^2)$
both expectation being with respect to the $N(x, \sigma^2)$
distribution. Therefore we simply sample
$\theta_1, \theta_2, \ldots$ from $N(x, \sigma^2)$ and use

$$\widehat{E^{\pi}(\theta \mid x)} = \frac{\sum_{i=1}^{m} \theta_i \left(\tau^2 + (\theta_i - \mu)^2\right)^{-1}}{\sum_{i=1}^{m} \left(\tau^2 + (\theta_i - \mu)^2\right)^{-1}}$$

as our Monte carlo estimate of $E^{\pi}(\theta \mid x)$.

But the problem has not been completely solved.

The sample of $\theta$'s generated from $N(x, \sigma^2)$
will tend to concentrate around $x$, whereas to
satisfactorily account for the contribution of the
Cauchy prior to the posterior mean, a significant
portion of the $\theta$'s should come from the tails of
the posterior distribution.

It may therefore appear that it is perhaps better $\wedge$ (6)
to express the posterior mean in the form

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

then sample $\theta$'s from Cauchy $(\mu, \tau)$ and use the approximation

$$\widehat{E^{\pi}(\theta|x)} = \frac{\sum_{i=1}^{m} \theta_i \, \exp\left(-\frac{(\theta_i-x)^2}{2\sigma^2}\right)}{\sum_{i=1}^{m} \exp\left(-\frac{(\theta_i-x)^2}{2\sigma^2}\right)}$$

However, this is also not totally satisfactory because the tails of the posterior distribution are not as heavy as those of the Cauchy prior, and hence there will be excess sampling from the tails relative to the center. To implication is that the convergence of the approximation is slower and hence a larger error in approximation (for a fixed $m$). Ideally, therefore, sampling should be from the posterior distribution itself for a satisfactory approximation. With this view in mind, a variation of the above theme has been develop. This is called the Monte Carlo importance sampling.

Consider $\quad E_f(h(\underline{x})) = \int\limits_{\chi} h(\underline{x}) f(\underline{x}) d\underline{x} \quad \cdots \cdots \quad (*)$

$= \cancel{\int f h f x}$

Suppose that it is difficult or expensive to sample directly from $f$, but there exists a probability density $u$ that is very close to $f$ from writ which it is easy to sample. Then can rewrite $(*)$ as

$$E_f(h(\underline{x})) = \int\limits_{\chi} h(\underline{x}) f(\underline{x}) d\underline{x}$$

$$= \int\limits_{\chi} h(\underline{x}) \frac{f(\underline{x})}{u(\underline{x})} u(\underline{x}) d\underline{x}$$

$$= \int\limits_{\chi} \{ h(\underline{x}) w(\underline{x}) \} u(\underline{x}) d\underline{x}$$

$$= E_u \{ h(\underline{x}) w(\underline{x}) \}$$

where $\quad w(\underline{x}) = \dfrac{f(\underline{x})}{u(\underline{x})}$ . Now generate i.i.d. observations $\underline{x}_1, \underline{x}_2 \cdots$ from the density $u$ and compute $\quad \overline{h w}_m = \dfrac{1}{m} \sum\limits_{i=1}^{m} h(x_i) w(x_i)$

The sampling density $u$ called the importance function.

**Example** Suppose $x_1, x_2, \cdots x_n$ are i.i.d $N(\theta, \sigma^2)$, where $\theta$ and $\sigma^2$ are unknown. Independent priors are assumed for $\theta$ and $\sigma^2$, where $\theta$ has a double exponential distribution with density $\exp(-|\theta|)/2$ and $\sigma^2$ has the prior density of $(1+\sigma^2)^{-2}$.

Neither of these is a standard prior, but robust choice of proper prior all the same. If the posterior mean of $\theta$ is of interest, then it is necessary to compute,

$$E^{\pi}(\theta|x) = \int_{-\infty}^{\infty} \int_{0}^{\infty} \theta \pi(\theta, \sigma^2|\underline{x}) \, d\theta \, d\sigma^2$$

Because $\pi(\theta, \sigma^2|\underline{x})$ is not a standard density, let us look for a standard density close to it.

Letting $\bar{x}$ denote the mean of the sample $x_1, x_2, \cdots, x_n$ and $S_n^2 = \sum_{i=1}^{n}(x_i - \bar{x})^2 / n$,

note that, $\pi(\theta, \sigma^2|\underline{x}) \propto (\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{\sigma}{2\sigma^2}\{(\theta-\bar{x})^2 + S_n^2\}\right) \times$

$$\exp(-|\theta|)(1+\sigma^2)^{-2}$$

$$= \left[S_n^2 + (\theta-\bar{x})^2\right]^{\frac{n}{2}+1} (\sigma^2)^{-(\frac{n}{2}+2)} \exp\left(-\frac{n}{2\sigma^2}\{(\theta-\bar{x})^2 + S_n^2\}\right)$$

$$\times \left\{[S_n^2 + (\theta-\bar{x})^2]^{-(\frac{n}{2}+1)} \exp(-|\theta|) \left(\frac{\sigma^2}{1+\sigma^2}\right)^2\right\}$$

$$\therefore \pi(\theta, \sigma^2|\underline{x}) \propto u_1(\sigma^2|\theta) \, u_2(\theta) \exp(-|\theta|) \left(\frac{\sigma^2}{1+\sigma^2}\right)^2$$

where $u_1(\sigma^2|\theta)$ is the density of inverse Gamma

with shape parameter $\frac{n}{2}+1$ and scale parameter $\frac{n}{2}\{(\theta-\bar{x})^2 + S_n^2\}$ and $u_2$ is the student's t density with d.f. $n+1$, location $\bar{x}$ and scale a multiple of $S_n$. It may be noted that the tails of $\exp(-|\theta|)\left(\frac{\sigma^2}{1+\sigma^2}\right)^2$ do not have much of an influence in the presence of $u_1(\sigma^2|\theta)\,u_2(\theta)$. Therefore,

$u(\theta,\sigma^2) = u_1(\sigma^2|\theta)\,u_2(\theta)$ may be chosen as a suitable importance function. This involves sampling $\theta$ first from the density $u_2(\theta)$ and given this $\theta$, sampling $\sigma^2$ from $u_1(\sigma^2|\theta)$. This is repeated to generate further value of $(\theta,\sigma^2)$. Finally, after generating $m$ of these pairs $(\theta,\sigma^2)$, the required posterior mean of $\theta$ is approximated by

$$E^{\pi}(\theta\,|\,\bar{x}) = \frac{\sum\limits_{i=1}^{m} \theta_i\, w(\theta_i,\sigma_i^2)}{\sum\limits_{i=1}^{m} w(\theta_i,\sigma_i^2)}.$$

where, $w(\theta,\sigma^2) = f(\bar{x}\,|\,\theta,\sigma^2)\,\pi(\theta,\sigma^2)\big/ u(\theta,\sigma^2)$

In some high-dimensional problems, a combination of numerical integration, Laplace approximation and Monte Carlo sampling seem to give appealing results.

# Markov Chain Monte Carlo Methods

A severe drawback of the standard Monte Carlo sampling or Monte Carlo importance sampling is that complete determination of the functional form of the posterior density is needed for their implimentation.

Situations where posterior distributions are imcompletely specified or are specified indirectly cannot be handled. One such instance is where the joint posterior distribution of the vector of parameters is specified in terms of several conditional and marginal distributions, but not directly.

This actually covers a very large range of Bayesian analysis because a lot of Bayesian modeling is hierarchical so that the joint posterior is difficult to calculate but the conditional posteriors given parameters at different levels of hierarchy are easier to write down and hence sample from.

For instance, consider the normal - Cauchy problem (example!). As shows This problem can be given a hierarchical structure wherein we have the normal model, the conjugate normal prior in the first stage with a hyperparameter for its variance and this hyperparameter again has the conjugate prior.

Similarly, consider example 2 where we have independent observations $X_i \sim$ Poisson $(\theta_i)$. Now suppose the prior on the $\theta_i$'s is a conjugate mixture. We again see that a hierarchical prior structure can lead to analytically tractable conditional posterior. It turns out that it is indeed possible in such cases to adopt an iterative Monte Carlo sampling scheme, which at the point of convergence will guarantee a random draw from the target joint posterior distribution. These iterative Monte Carlo procedures typically generate a random sequence with the Markov property such that this Markov chain is ergodic with the limiting distribution being the target posterior distribution. This is actually a whole class of such iterative procedures collectively called Markov chain Monte Carlo (MCMC) ~~produ~~ procedures. Different procedure from this class are suitable for different situations.

As mentioned above, convergence of a random sequence with the Markov property is being utilized in this procedure, and hence some basic understanding of Markov chains is required.

# Law of large numbers for Markov chains

Let $\{X_n\}_{n \geq 0}$ be a Markov chain with a countable state space $S$ and a transition probability matrix $P$. Further, suppose it is irreducible and has a stationary probability distribution $\pi = (\pi_i : i \in S)$. Then, for any bounded function $h : S \longrightarrow R$ and for any initial distribution of $X_0$

$$\frac{1}{n} \sum_{i=0}^{n-1} h(X_i) \longrightarrow \sum_{j} h(j) \pi_j$$

in probability as $n \longrightarrow \infty$.

A similar law of large numbers (LLN) holds when the state space $S$ is not countable. The limit value will be the integral of $h$ with respect to the stationary distribution $\pi$.

A sufficient condition for the validity of this LLN is that the Markov chain $\{X_n\}$ be Harris irreducible and have a stationary distribution $\pi$. Harris irreducibility is the notion of irreducibility in general case.

To see how this is useful, consider the following:
Given a probability distribution $\pi$ on a set $S$, and a function $h$ on $S$, suppose it is desired to compute the "integral of $h$ with respect to $\pi$", which reduces to $\sum_j h(j) \pi_j$ in the countable case. Look for an irreducible Markov chain $\{X_n\}$ with state space $S$ and stationary distribution $\pi$. Then, starting from some initial value $X_0$, run the Markov chain $\{X_j\}$ for a period of time, say, $0, 1, 2, \ldots n-1$ and consider as an estimate

$$\mu_n = \frac{1}{n} \sum_0^{n-1} h(x_j).$$

By LLN, this estimate $\mu_n$ will be close to $\sum h(j) \pi_j$ for large $n$.

This technique is called Markov chain Monte Carlo (MCMC). For example, if one is interested in $\pi(A) = \sum_{j \in A} \pi_j$ for some $A \subseteq S$ then by LLN this reduces to $\pi_n(A) = \frac{1}{n} \sum_0^{n-1} I_A(x_j) \to \pi(A)$ in probability as $n \to \infty$, ~~for any initial distribution~~ of ~~$X_0$~~. Where $I_A(x_j) = 1$ if $x_j \in A$ and $0$ otherwise.

An irreducible Markov chain $\{X_n\}$ with a countable state space $S$ is called aperiodic if for some $i \in S$ the g.c.d $\{ n : p_{ii}^{(n)} > 0 \} = 1$. Then, in addition to the LLN, the following result on the convergence of $P(X_n = j)$ holds.

$$\sum_{j} | P(X_n = j) - \pi_j | \to 0 \qquad \text{as } n \to \infty,$$

for any initial distribution of $X_0$.

In the other words, for large $n$ the probability distribution of $X_n$ will be close to $\pi$. There exists a result similar to above for the general state space case that asserts that under suitable conditions, the probability distribution of $X_n$ will be close to $\pi$ as $n \to \infty$.

This suggests that instead of doing over run of length $n$, one could do $N$ dependent runs each of length $m$ so that $n = Nm$ and then from the $i^{th}$ run use only the $m^{th}$ observation, say, $X_{m,i}$, and consider the estimate

$$\tilde{\mu}_{N,m} = \frac{1}{N} \sum_{i=1}^{N} h(X_{m,i})$$

Other variations exist as well.

# Metropolis - Hasting Algorithm

Here, we discuss a very general MCMC method with wide applications. The idea here is not to directly simulate from the 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 probability distribution on $S$. We shall call $\pi$ the target distribution. Let $Q \equiv ((q_{ij}))$ be a transition probability matrix such that for each $i$, it is computationally easy to generate a sample from the distribution $\{q_{ij} : j \in S\}$. Let us generate a Markov chain $\{X_n\}$ as follows. If $X_n = i$, first sample from the distribution $\{q_{ij} : j \in S\}$ and denote that observation $Y_n$. Then, choose $X_{n+1}$ from the two values $X_n$ and $Y_n$ according to

$$P(X_{n+1} = Y_n \mid X_n, Y_n) = \rho(X_n, Y_n)$$

$$P(X_{n+1} = X_n \mid X_n, Y_n) = 1 - \rho(X_n, Y_n)$$

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

$$\rho(i,j) = \min\left\{ \frac{\pi_j \, q_{ji}}{\pi_i \, q_{ij}} , 1\right\} \quad \text{for all } (i,j) \text{ such that } \pi_i q_{ij} > 0.$$

Note that $\{X_n\}$ is a Markov chain with transition probability matrix $P = ((P_{ij}))$ given by

$$p_{ij} = \begin{cases} q_{ij} \, \rho_{ij} & j \neq i \\ \\ 1 - \sum_{k \neq i} p_{ik} & j = i \end{cases}$$

Q is called the "proposal transition probability" and $\rho$ the "acceptance probability".

A significant feature of this transition mechanism P is that P and $\pi$ satisfy

$$\pi_i \, p_{ij} = \pi_j \, p_{ji} \qquad \text{for } i, j \, .$$

This implies that for any $j$,

$$\sum_i \pi_i \, p_{ij} = \pi_j \sum_i p_{ji} = \pi_j$$

or, $\pi$ is a stationary probability distribution for P.

Now assume that S is irreducible with respect to Q and $\pi_i > 0$ for all $i$ in S. It can be then be shown that P is irreducible, and because it has a stationary distribution $\pi$, LLN is available. This algorithm is thus a very flexible and useful one. The choice of Q is subject only to the condition that S is irreducible with respect to Q. Clearly, it is no loss of generality to assume that $\pi_i > 0$ for all $i$ in S.

A sufficient condition for the aperiodicity of $P$ is that $p_{ii} > 0$ for some $i$ or equivalently

$$\sum_{j \neq i} q_{ij} \, \rho_{ij} < 1.$$

A sufficient condition for this is that there exists a pair $(i, j)$ such that $\pi_i \, q_{ij} > 0$ and

$$\pi_j \, q_{ji} < \pi_i \, q_{ij} \, .$$

Recall that if $P$ is aperiodic, then both the LLN and $\sum_j |P(x_n = j) - \pi_j| \to 0$ as $n \to \infty$ hold. If $S$ is not finite or countable but is a continum and the target distribution $\pi(\cdot)$ has a density $p(\cdot)$, then one proceeds as follows: Let $Q$ be a transition function such that for each $x$, $Q(x, \cdot)$ has a density $q(x, y)$. Then proceed as in the discrete case but set the acceptance probability $\rho(x, y)$ to be $\rho(x, y) = \min \left\{ \dfrac{p(y) \, q(y, x)}{p(x) \, q(x, y)}, 1 \right\}$ for all $(x, y)$ such that $p(x) \, q(x, y) > 0$.

A particularly useful feature of the above algorithm is that it is enough to know $p(\cdot)$ upto a multiplicative constant. This assures us that in Bayesian applications it is not necessary to have the normalizing constant of the posterior density.

# Gibbs Sampling

Most of the new problems that Bayesians are asked to solve are high-dimensional. Application areas such as micro-arrays and image processing are some examples. Bayesian analysis of such problems invariably involve target (posterior) distributions that are high-dimensional multivariate distributions. In image processing, for example, typically one has $N \times N$ grid of pixels with $N = 256$ and each pixel has $k \geq 2$ possible values. Thus each configuration has $(256)^2$ components and the state space $S$ has $k^{(256)^2}$ configurations. To simulate a random configuration from a target distribution over such a large $S$ is not an easy task.

The Gibbs sampler is a technique especially suitable for generating an irreducible aperiodic Markov chain that has as its stationary distribution a target distribution in a high-dimensional space but having some special structure. The most interesting aspect of this technique is that to run this Markov chain, it suffices to generate observations from univariate distributions.

The Gibbs sampler in the context of a bivariate probability distribution can be described as follows. Let $\pi$ be a target probability distribution of a bivariate random vector $(x, y)$. For each $x$, let $P(x, \cdot)$ be the conditional probability distribution of $Y$ given $X = x$. Similarly, let $Q(y, \cdot)$ be the distribution $X$ given $Y = y$.

Note that for each $x$, $P(x, \cdot)$ is a univariate distribution, and for each $y$, $Q(y, \cdot)$ is also a univariate distribution. Now generate a bivariate Markov chain $Z_n = (X_n, Y_n)$ as follows:

Start with some $X = x_0$. Generate an observation $Y_0$ from the distribution $P(x_0, \cdot)$. Then generate an observation $x_1$ from $Q(Y_0, \cdot)$. Next generate an $Y_1$ from $P(x_1, \cdot)$ and so on.

If $\pi$ is a discrete distribution concentrated on $\{(x_i, y_i) : 1 \le i \le k, 1 \le j \le L\}$ and if

$$\pi_{ij} = \pi(x_i, y_j) \quad \text{then} \quad P(x_i, y_j) = \frac{\pi_{ij}}{\pi_{i\cdot}} \quad \text{and}$$

$$Q(y_j, x_i) = \frac{\pi_{ij}}{\pi_{\cdot j}},$$

Where $\pi_{i\cdot} = \sum_j \pi_{ij}$ and $\pi_{\cdot j} = \sum_i \pi_{ij}$

Thus the transition probability matrix

$$R = ((r_{(ij),(kl)})) \text{ for the } \{Z_n\} \text{ chain is given}$$

by 

$$r_{(ij)(kl)} = Q(y_j, x_k) P(x_k, y_l)$$

$$= \frac{\pi_{kj}}{\pi_{\cdot j}} \frac{\pi_{kl}}{\pi_{k\cdot}}$$

It can be verified that this chain is irreducible, aperiodic, and has $\pi$ as its stationary distribution. Thus LLN holds in this case.. Thus for large $n$, $Z_n$ can be viewed as a sample from a distribution that is close to $\pi$ and one can approximate $\sum_{i,j} h(i,j) \pi_{ij}$ by $\sum_{i=1}^{n} h(x_i, y_i)/n$

<u>Example</u>    Consider sampling from $\binom{x}{z} \sim N_2\left(\binom{0}{0}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$. Note that the conditional distribution of $X$ given $Y=y$ and that of $Y$ given $X=x$ are

$$X|Y=y \sim N(\rho y, 1-\rho^2) \text{ and } Y|X=x \sim N(\rho x, 1-\rho^2)$$

Using this property, Gibbs sampling proceeds as described below to generate $(X_n, Y_n)$, $n = 0, 1, 2 \ldots$ by starting from an arbitrary value $x_0$ for $X_0$ and repeating the following steps for $i = 0, 1, 2 \ldots n$.

1. Given $x_i$ for $X$, draw a random deviate from $N(\rho x_i, 1-\rho^2)$ and denote it by $Y_i$.

2. Given $y_i$ for $Y$, draw a random deviate from $N(\rho y_i, 1-\rho^2)$ and denote it by $X_{i+1}$.

The theory of Gibbs sampling tells us that if $n$ is large, then $(x_n, y_n)$ is a random draw from a distribution that is close to $N_2\left(\left\{\begin{smallmatrix}0\\0\end{smallmatrix}\right\}, \left[\begin{smallmatrix}1&\rho\\\rho&1\end{smallmatrix}\right]\right)$.

To see why Gibbs sampler works here, recall that a sufficient condition for the LLN and the limit result is that an appropriate irreducibility condition holds and a stationary distribution exists. From steps 1 and 2 above and using the conditional distributions $X/Y = y$ and $Y | X = x$, one has

$$Y_i = \rho x_i + \sqrt{1-\rho^2} \cdot \eta_i$$

and

$$X_{i+1} = \rho Y_i + \sqrt{1-\rho^2} \cdot \xi_i ,$$

Where $\eta_i$ and $\xi_i$ are independent standard normal random variables independent of $X_i$. Thus the sequence $\{X_i\}$ satisfies the stochastic difference equation $X_{i+1} = \rho^2 X_i + U_{i+1}$,

Where $U_{i+1} = \rho \sqrt{1-\rho^2} \eta_i + \sqrt{1-\rho^2} \xi_i$.

Because $\eta_i, \xi_i$ are independent $N(0,1)$ random variables, $U_{i+1}$ is also a normally distributed random variable with mean 0 and variance $\rho^2(1-\rho^2) + (1-\rho^2) = 1-\rho^4$. Also $\{U_i\}_{i \geq 1}$ being i.i.d. makes $\{X_i\}_{i \geq 0}$ a Markov chain. It turns out that the irreducibility condition holds here. Also the standard $N(0,1)$ distribution a stationary distribution for $\{X_n\}$.