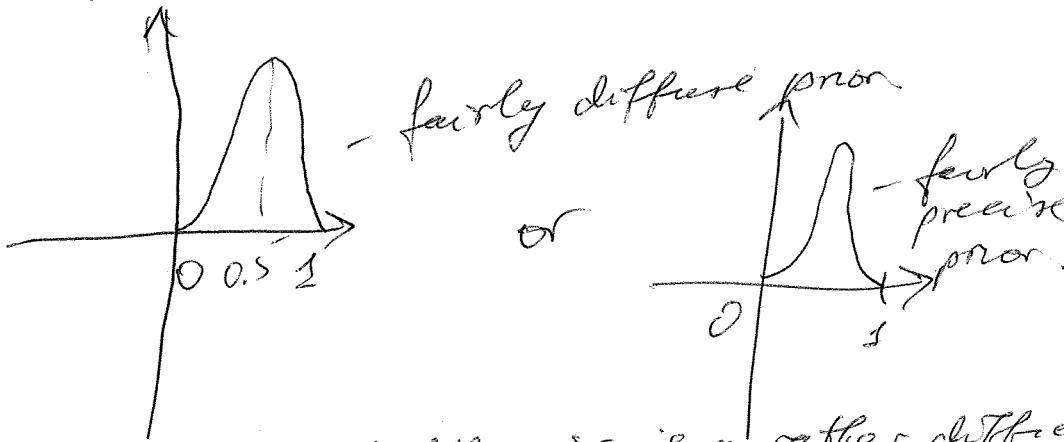


7.1 Prior and Posterior Distributions

Model $\{f_{\theta}: \theta \in \mathcal{S}\}$ for the data $s \in S$
and the prior probability measure π for θ .

Ex. If $\mathcal{S} = [0, 1]$, θ is the prob. of success
on a coin toss, we may have a prior like



Assignment of the prior is a rather difficult problem. Thus, the full set of ingredients in the Bayesian approach are: 1) a distribution for θ - prior π
2) a set of conditional distributions for the data given θ ,
i.e. $\{f_{\theta}(s|\theta)\}$. \Rightarrow the joint probability distn for (s, θ) is $\pi(\theta)f_{\theta}(s)$. If the prior distribution is absolutely continuous, the marginal dist. for s is $m(s) = \int \pi(\theta)f_{\theta}(s)d\theta$ - this is the prior predictive distribution of s the data.
Next, posterior dist. of θ is described as

$$\pi(\theta|s) = \frac{\pi(\theta)f_{\theta}(s)}{m(s)}, \text{ this is an axiom based on the}$$

conditional probability principle, not a theorem!
 $m(s)$ is the inverse normalizing constant for the posterior density since $\pi(\theta|s) \propto \pi(\theta)f_{\theta}(s)$. In many examples, $m(s)$ need not be computed because we recognize the functional form as a function of θ .

Ex. 7.1.1. Bernoulli model -

Let $(x_1, x_2, \dots, x_n) \sim \text{Ber}(\theta)$, $\theta \in [0, 1]$.

Let $\tau = \text{Beta}(\alpha, \beta)$; then, the posterior of θ is proportional to the likelihood $\prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\bar{x}} (1-\theta)^{n-\bar{x}}$

multiplied by the prior $B^{-1}(\alpha, \beta) \propto (\alpha + \beta)^{\alpha + \beta - 2}$

⇒ The product is proportional to

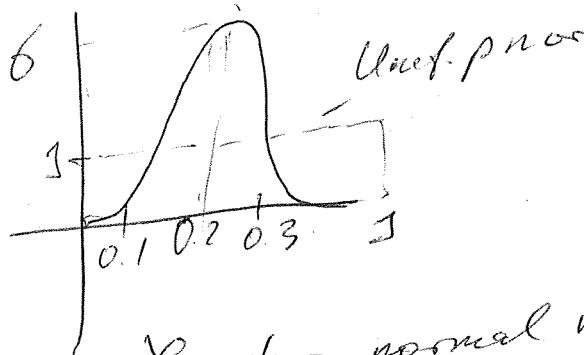
$$\theta^{n\bar{x} + \alpha - 1} (1-\theta)^{n(1-\bar{x}) + \beta - 1}$$

- Beta($n\bar{x} + \alpha, n(1-\bar{x}) + \beta$), no need to compute

$m(x_1, \dots, x_n)$. Let us say, we observe $n\bar{x} = 10$ in a sample

of $n=40$, and $\alpha = \beta = 1 \Rightarrow$ uniform prior on θ . Then,

the posterior of θ is given by Beta(11, 31). $\downarrow \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$



from here comes
Beta(α, β) distribution

Ex. 7.1.2. Location normal model

$(x_1, \dots, x_n) \sim N(\mu, \sigma^2)$, σ^2 is known

$$L(\mu | x_1, \dots, x_n) = \exp\left(-\frac{n}{2\sigma^2}(\bar{x} - \mu)^2\right)$$

Let the prior μ be $N(\mu_0, \tau_0^{-2})$ for some specified μ_0 and

τ_0^2 . The posterior density of the μ is then

$$-\frac{1}{2\tau_0^2}(\mu - \mu_0)^2 - \frac{1}{2\sigma^2}(\bar{x} - \mu)^2 = e^{-\frac{\mu_0^2}{2\tau_0^2} - \frac{n\bar{x}^2}{2\sigma^2}}$$

$$e^{-\frac{1}{2}(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2})[\mu^2 - 2\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2}\bar{x}\right)\mu]} =$$

$$\times e^{-\frac{1}{2}\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)\left(\mu - \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2}\bar{x}\right)\right)^2}$$

$$= e^{-\frac{1}{2}\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2}\bar{x}\right)^2} \times e^{-\frac{1}{2}\left(\frac{\mu_0^2}{\tau_0^2} + \frac{n\bar{x}^2}{\sigma^2}\right)}$$

$$\therefore \text{This is } \propto N\left(\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2}\bar{x}\right), \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right)$$

Posterior mean is \bar{x} weighted average of prior mean μ_0 and the sample mean \bar{x}

$$\text{weights } \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1} \frac{1}{\tau_0^2} \text{ and } \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma_0^2}\right)^{-1} \frac{n}{\sigma_0^2}$$

Also, the posterior variance is smaller than the variance of the sample mean. The more diffuse the prior is, (\Rightarrow the larger τ_0^2 is) - the less influence the prior has.

E.g. if $n=20$, $\sigma_0^2 = 1$, $\tau_0^2 = 1$
the ratio of the posterior variance to the sample mean variance is $\frac{20}{21} \approx 0.95$

Ex 7.13 Multinomial model.

A categorical response x that takes k values

$\forall i \in S = \{1, 2, \dots, k\}$. Say, a bowl contains chips labelled $1, 2, \dots, k$. A proportion θ_i of the chips are labelled i . We randomly draw a chip, observing its label. With unknown θ_i , we have

$$\{P(\theta_1, \dots, \theta_k) : (\theta_1, \dots, \theta_k) \in \Omega\} \text{ where } P(\theta_1, \dots, \theta_k) = \prod_{i=1}^k \theta_i^{x_i}$$

$$\text{and } \Omega = \{(\theta_1, \dots, \theta_k) : 0 \leq \theta_i \leq 1, i=1, \dots, k, \sum_{i=1}^k \theta_i = 1\}$$

A sample (s_1, \dots, s_n) is observed; let the i^{th} category count be x_i , $\Rightarrow L(\theta_1, \dots, \theta_k | s_1, \dots, s_n) = \theta_1^{x_1} \theta_2^{x_2} \dots \theta_k^{x_k}$

The prior: $(\theta_1, \dots, \theta_{k-1}) \sim \text{Dirichlet}(\alpha_1, \alpha_2, \dots, \alpha_k)$

$$\text{with density } \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} \frac{\alpha_1^{\alpha_1-1} \alpha_2^{\alpha_2-1} \dots \alpha_k^{\alpha_k-1}}{\theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k}}$$

$\alpha_i > 0$ are constants that represent the prior belief about $(\theta_1, \dots, \theta_{k-1})$. The choice $\alpha_1 = \alpha_2 = \dots = \alpha_k = 1$ - uniform distribution

\Rightarrow The posterior density of $(\theta_1, \dots, \theta_{k-1})$ is the

$$\propto \theta_1^{\alpha_1+x_1-1} \theta_2^{\alpha_2+x_2-1} \dots \theta_k^{\alpha_k+x_k-1} \sim \text{Dirichlet}(\alpha_1+x_1, \dots, \alpha_k+x_k)$$

Location-Scale Normal Model. -4-

Now, $(x_1, x_n) \sim N(\mu, \sigma^2)$ $\frac{n}{2\sigma^2} (\bar{x} - \mu)^2 = \frac{n-1}{2\sigma^2} s^2$

 $L(\mu, \sigma^2 | x_1, x_n) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} (\bar{x} - \mu)^2}$

1. $\mu/\sigma^2 \sim N(\mu_0, \frac{1}{\tau_0^2} \sigma^2)$,
2. $\frac{1}{\sigma^2} \sim \text{Gamma}(\alpha_0, \beta_0)$ - inverse Gamma -

From here (no details):

 $\mu(\sigma^2 | x_1, x_n) \sim N(\mu_x, \left(n + \frac{1}{\tau_0^2}\right)^{-1} \sigma^2)$
 $\text{and } \frac{1}{\sigma^2} | x_1, x_n \sim \text{Gamma}(\alpha_0 + n/2, \beta_x)$

Where $\mu_x = \left(n + \frac{1}{\tau_0^2}\right)^{-1} \left(\frac{\mu_0}{\tau_0^2} + n\bar{x}\right)$

and $\beta_x = \beta_0 + \frac{n-1}{2} s^2 + \frac{1}{2} \frac{n(\bar{x} - \mu_0)^2}{1 + n\tau_0^2}$

As $\tau_0 \rightarrow \infty$ (the prior is increasingly diffuse)

the conditional posterior distribution of μ given σ^2
 $\xrightarrow{D} N(\bar{x}, \frac{\sigma^2}{n})$ because $\mu_x \rightarrow \bar{x}$ and $\left(n + \frac{1}{\tau_0^2}\right)^{-1} \rightarrow \frac{1}{n}$

Furthermore, as $\tau_0 \rightarrow \infty$ and $\beta_0 \rightarrow 0$, the marginal posterior of $\frac{1}{\sigma^2} \xrightarrow{D} \text{Gamma}(\alpha_0 + \frac{n}{2}, \frac{(n-1)s^2}{2})$ because
 $\beta_x \rightarrow (n-1)s^2/2$. Note that there's no proper prior distribution as $\tau_0 \rightarrow \infty$ and $\beta_0 \rightarrow 0$

Corresponds to the diffuse prior
 Where for $\text{Gamma}(\alpha_0, \beta_0)$ $\text{Var} = \frac{\alpha_0}{\beta_0^2} \xrightarrow{\beta_0 \rightarrow 0} \infty$

7.2. Inferences based on the posterior distribution.

7.2.1. Estimation

To estimate $\bar{\theta}(\Theta)$, two most common approaches are: 1) to compute the posterior density of $\bar{\theta}(\Theta)$ and use the posterior mode. To do so, one needs to maximize $w(\bar{\theta}|S)$, or equivalently $m(S)w(\bar{\theta}|S)$ (so as not to compute the normalizing constant), as a function of $\bar{\theta}$. Moreover, in general, any 1-1 increasing function of $\bar{\theta}(S)$ can be maximized to yield the same result.

2) Another common alternative is to use the posterior mean $E(\bar{\theta}(\Theta)|S)$. When the posterior distribution is symmetric and the expectation is finite, this is the same as the mode.

Ex. 7.2.2. Bernoulli model -

$X_1, \dots, X_n \sim \text{Ber}(\Theta)$, $\Theta \in [0, 1]$ unknown; $\Theta \sim \text{Beta}(\alpha, \beta)$

$$\text{The posterior distribution of } \Theta \text{ is } \text{Beta}(n\bar{x} + \alpha, n(1-\bar{x}) + \beta).$$

$$\text{Then } E(\Theta|X_1, \dots, X_n) = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(n\bar{x}+\alpha)\Gamma(n(1-\bar{x})+\beta)} \int_0^1 \theta^{n\bar{x}+\alpha} (1-\theta)^{n(1-\bar{x})+\beta} d\theta =$$

$$= \frac{\Gamma(n+\alpha+\beta)}{\Gamma(n\bar{x}+\alpha)\Gamma(n(1-\bar{x})+\beta)} \frac{\Gamma(n\bar{x}+\alpha+1)\Gamma(n(1-\bar{x})+\beta+1)}{\Gamma(n+\alpha+\beta+2)} =$$

$$= \frac{n\bar{x} + \alpha}{n + \alpha + \beta}; \text{ when } \alpha = \beta = 1 \text{ (uniform prior), the}$$

$$\text{posterior expectation is } E(\Theta|X_1, \dots, X_n) = \frac{n\bar{x}}{n+2}.$$

To determine the posterior mode, we have to maximize $\ln \Theta^{n\bar{x}+\alpha-1} (1-\Theta)^{n(1-\bar{x})+\beta-1}$; direct approach yields

$$\hat{\Theta} = \frac{n\bar{x} + \alpha - 1}{n + \alpha + \beta - 2}, \text{ which is negative whenever } \alpha > 1, \beta < 1.$$

This last restriction implies that the mode of the prior is in $(0, 1)$ and not at 0 or 1. When $\alpha = \beta = 1$, the mode is $\hat{\Theta} = \bar{x}$ - the same as MLE. When n is large, both are approx. the same and close to the MLE.

7.8.3 Location normal model
 Now, $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, μ is unknown, σ^2 is known.
 The prior on mean μ is $N(\mu_0, \tau_0^2)$

The posterior distribution of μ is

$$N\left(\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1} \left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2} \bar{x}\right), \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1}\right)$$

This is symmetric about its mode, so the posterior mode and the mean are the same: $\left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right)^{-1} \left(\frac{\mu_0}{\tau_0^2} + \frac{n}{\sigma^2} \bar{x}\right) =$

- the weighted average of the prior mean and the sample mean.

Remarks: 1) When $n \rightarrow \infty$ we get \bar{x} - MLE
 2) When $\tau_0^2 \rightarrow \infty$ (diffuse prior),

get \bar{x} - MLE again
 3) The ratio of the sampling variance of \bar{x} to the posterior variance of μ is $\frac{\sigma^2}{n} \left(\frac{1}{\tau_0^2} + \frac{n}{\sigma^2}\right) = 1 + \frac{\sigma^2}{n\tau_0^2} > 1$.

The closer τ_0^2 is to 0, the larger this ratio is
 The closer τ_0^2 is to 0, posterior mean $\rightarrow \mu_0$

Numerical example: Average weight μ of adult cows in a large herd. Farmer weighs 10 randomly chosen cows from the herd; $\bar{Y} = 410$ kg / 925 lbs; let $\sigma^2 = 20$ kg and the prior distribution is normal $\Rightarrow \bar{Y}$ = MLE. Prior belief of the farmer: the weight should be ≈ 420 kg and unlikely to differ from it by more than 30 kg \rightarrow Prior $N(420, 10^2)$ due to the empirical three Sigma rule; see posterior Bayes's

Then $\frac{100}{400/10 + 100} 410 + \frac{400/10}{400/10 + 100} 420 = 412.9$ kg

\rightarrow the posterior mean is $\frac{\mu_0}{\tau_0^2} \frac{\tau_0^2}{\sigma^2 + n\tau_0^2} + \frac{n}{\sigma^2} \frac{\bar{x}}{\sigma^2 + n\tau_0^2}$; here $\tau_0^2 = 10^2 = 100$

$$= \mu_0 + \frac{\sigma^2/n}{\sigma^2/n + \tau_0^2} + \frac{\tau_0^2}{\sigma^2/n + \tau_0^2} \bar{x} ; \quad \begin{aligned} \sigma^2 &= 20^2 = 400 \\ n &= 10 \end{aligned}$$