

Lecture 8A

MTH6102: Bayesian Statistical Methods

Eftychia Solea

Queen Mary University of London

2023

Today's agenda

Today's lecture will

- Review of noninformative priors
- Learn informative priors
- Be able to make a reasonable choice of informative prior, based on external data.
- Learn that the choice of prior affects the posterior.
- See that more data lessens the dependence of the posterior on the prior.

- The prior distribution plays a defining role in Bayesian analysis.
- There are two types of priors: **noninformative** and **informative**.

Noninformative prior

- A **noninformative prior** represents our ignorance or lack of information about θ before the experiment.
- Non-informative prior: “let the data speak for themselves”.
- An obvious candidate for a noninformative prior is to use a **flat or uniform prior** over some range, $p(\theta) \propto c$.
- It is flat relative to the range of the likelihood. Assumes that every hypothesis is equally probable.
- **Flat priors are not invariant under nonlinear one-to-one transformations g**

Noninformative Jeffreys prior

- Another example of noninformative prior is **Jeffreys prior** defined as

$$p_J(\theta) \propto \sqrt{I(\theta)},$$

where $I(\theta)$ is the Fisher information function given by (under some regularity conditions)

$$I(\theta) = -E\left[\frac{d^2}{d\theta^2} \log p(X|\theta)\right],$$

and $p(X|\theta)$ is the likelihood function.

$\cdot X \sim \text{binomial}(n, \theta), \quad p_J(\theta) \sim \text{beta}\left(\frac{1}{2}, \frac{1}{2}\right)$

Noninformative Jeffreys prior

- Jeffreys prior is invariant under nonlinear, smooth, one-to-one transformations g because

$$I(\psi) = I(\theta) \left(\frac{d\theta}{d\psi} \right)^2$$

invariance
property

where $\psi = g(\theta)$.

Noninformative prior

- **Advantages**

- ① sometimes used as benchmark that will not reflect the bias of the analyst.
- ② appropriate when little is known on the underlying processes.
- ③ can be used in situations where scientific objectivity is at a premium, for example, when presenting results to a regulator or in a scientific journal,

- **Disadvantages**

- ① may lead to improper priors
- ② Flat priors are not invariant under nonlinear one-to-one transformations g
- ③ the definition of knowing little may depend on different parameterizations (should θ be assumed to be uniform or should perhaps the logarithm of θ be assumed to be uniform?)

Informative prior

- Informative priors include some judgement concerning plausible values of the parameters based on external information.
- Informative priors can be based on pure judgement, a mixture of data and judgement, or external data alone.
- An informative prior distribution is one in which the probability mass is concentrated in some subset of the possible range for the parameters.

- There are many ways to build an informative prior. For example, using summary statistics, published estimates, intervals or standard errors.
- We can match these quantities to the mean, median standard deviation or percentiles of the prior distribution.

Example: Building an informative prior

Exponential/Gamma example

- Let $t_1, \dots, t_n \sim \text{Exp}(\lambda)$ denote the lifetimes of lightbulbs.
- The gamma distribution is conjugate to the exponential likelihood for λ (failure rate).
- Suppose we have external information from other similar bulbs with observed failure rates r_1, \dots, r_K .
- Let m and u be the mean and variance of r_1, \dots, r_K , respectively.
- **Goal:** Build a prior gamma(α, β) distribution for λ using external information.

Example: Building an informative prior

Exponential/Gamma example

- We can use the method of moments to match the mean and the variance of the prior gamma distribution with the corresponding m and u

- That is

$$m = \frac{\alpha}{\beta}, \quad u = \frac{\alpha}{\beta^2}$$

- Solve for α and β

$$\beta = \frac{m}{u}, \quad \alpha = \frac{m^2}{u}$$

- Thus, our prior for λ is gamma($\frac{m^2}{u}$, $\frac{m}{u}$).

sample mean

Gamma mean

Gamma variance

m, u observed

$$m = \frac{a}{b} \Rightarrow a = mb$$

Then, $u = \frac{mb}{b^2}$

$$\Rightarrow b = \frac{m}{u}$$

$$a = \frac{m^2}{u}$$

- **Advantages**

- ① often analytically convenient (esp for conjugate priors).
- ② can take advantage of your informed understanding, beliefs, experience and external data

- **Disadvantages**

- ① not always easy to quantify the state of knowledge

Board example: Building an informative prior

Binomial/beta example

- Suppose we flip the coin n times and observe k heads with q the probability of heads.
- A beta(α, β) distribution is chosen as the prior distribution for q .
- Based on external information and published statistics, the prior mean is 0.4 and the prior standard deviation is 0.2.
- Find the prior distribution corresponding to this belief.
- See also **Question 2, final exam 2020**

Why?

prior mean $\frac{a}{a+b}$

prior variance $\frac{ab}{(a+b)^2(a+b+1)}$

We want to find $a > 0$ and $b > 0$ such that

$$\frac{a}{a+b} = \frac{4}{10} = \frac{2}{5} \quad (1)$$

$$\frac{ab}{(a+b)^2(a+b+1)} = \left(\frac{2}{10}\right)^2 = \frac{4}{100} = \frac{1}{25} \quad (2)$$

Solve for a and b .

From (1) $5a = 2a + 2b \Rightarrow 3a = 2b \Rightarrow a = \frac{2b}{3} \quad (3)$

Replace (3) into (2) to find,

$$\frac{\left(\frac{2b}{3}\right)b}{\left(\frac{2b}{3} + b\right)^2 \left(\frac{2b}{3} + b + 1\right)} = \frac{1}{25}$$

$$\Rightarrow \frac{\frac{2b^2}{3}}{25b^2 \left(\frac{5b}{3} + 1\right)} = \frac{1}{25}$$

$$\Rightarrow \frac{2 \cdot \cancel{b^2}}{\cancel{b^2} \cdot 25 \left(\frac{5b}{3} + 1\right)} = \frac{1}{25} \Rightarrow 6 = \frac{5b}{3} + 1$$
$$\Rightarrow \boxed{b=3} \quad \boxed{a=2}$$

Weakly informative prior distributions

- Instead of trying to make the prior completely uninformative, an alternative is to convey some information about the plausible range of the parameters, e.g., exclude implausible values.
- Otherwise let the data speak for themselves.
- For models with large numbers of parameters, adding a little prior information may help with numerical stability.

The choice of prior affects the posterior

- In the Bayesian framework, all our inferences about θ are based on the posterior distribution $p(\theta | y)$.

$$p(\theta | y) \propto p(\theta) p(y | \theta)$$

Posterior distribution \propto prior distribution \times likelihood

- Including summaries such as point estimates and credible intervals.
- So our inference depends on the prior distribution as well as the data via the likelihood.
- The choice of prior affects the posterior.
- More data, lessens the dependence of the posterior on the prior.

Normal example, known variance

- Observed data $y_1, \dots, y_n \sim N(\mu, \sigma^2)$.
- Prior distribution $\mu \sim N(\mu_0, \sigma_0^2)$.
- The posterior distribution is

$$\mu \mid y \sim N(\mu_1, \sigma_1^2)$$

$$\mu_1 = \left(\frac{\mu_0}{\sigma_0^2} + \frac{n\bar{y}}{\sigma^2} \right) / \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)$$

$$\sigma_1^2 = 1 / \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right)$$

Normal example: Posterior mean

- The posterior mean μ_1 can be written as a weighted average of the prior mean μ_0 and the sample mean \bar{y}

$$\mu_1 = (1 - w)\mu_0 + w\bar{y},$$

where

$$w = \frac{n}{\sigma^2} / \left(\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2} \right) = \frac{\sigma_0^2}{\frac{\sigma^2}{n} + \sigma_0^2}.$$

$w \rightarrow 1$ as $n \rightarrow \infty$ or $\sigma_0 \rightarrow \infty$, so the posterior mean approaches the sample mean.

$$n \rightarrow \infty \quad \mu_1 \rightarrow \bar{y}$$

$$\sigma_0 \rightarrow \infty \quad \mu_1 \rightarrow \bar{y}$$

Normal example: Likelihood and prior

- When deriving the posterior distribution, we saw that the likelihood

$$p(y | \mu)$$

is proportional to a

$$N\left(\bar{y}, \frac{\sigma^2}{n}\right) \text{ pdf for } \mu$$

if considered as a function of μ .

- If we compare $\frac{\sigma^2}{n}$ to the prior variance σ_0^2 , this helps to understand how the posterior behaves.

p(p|y) \propto likelihood \times prior

$$\propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right\} \exp\left\{-\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2\right\}$$

$$= \exp\left\{-\frac{1}{2\sigma^2} (n\mu^2 - 2n\bar{y}\mu + \sum y_i^2)\right\} \mathcal{N}(\mu_0, \sigma_0^2)$$

constant *function of μ*

$$\propto \exp\left\{-\frac{1}{2\sigma^2} (n\mu^2 - 2n\bar{y}\mu)\right\} \cdot \mathcal{N}(\mu_0, \sigma_0^2)$$

function of μ (complete the square)

$$= \exp\left\{-\frac{n}{2\sigma^2} (\mu^2 - 2\bar{y}\mu)\right\} \cdot \mathcal{N}(\mu_0, \sigma_0^2)$$

$$= \exp\left\{-\frac{n}{2\sigma^2} (\mu^2 - 2\bar{y}\mu + \bar{y}^2 - \bar{y}^2)\right\} \mathcal{N}(\mu_0, \sigma_0^2)$$

constant

$$\propto \exp\left\{-\frac{n}{2\sigma^2} (\underbrace{\mu^2 - 2\bar{y}\mu + \bar{y}^2}_{(\mu - \bar{y})^2})\right\} \mathcal{N}(\mu_0, \sigma_0^2)$$

$$= \exp\left\{-\frac{n}{2\sigma^2} (\mu - \bar{y})^2\right\} \mathcal{N}(\mu_0, \sigma_0^2)$$

$$= \exp\left\{-\frac{1}{2\sigma^2/n} (\mu - \bar{y})^2\right\} \mathcal{N}(\mu_0, \sigma_0^2)$$

Normal
likelihood
as a function
of μ

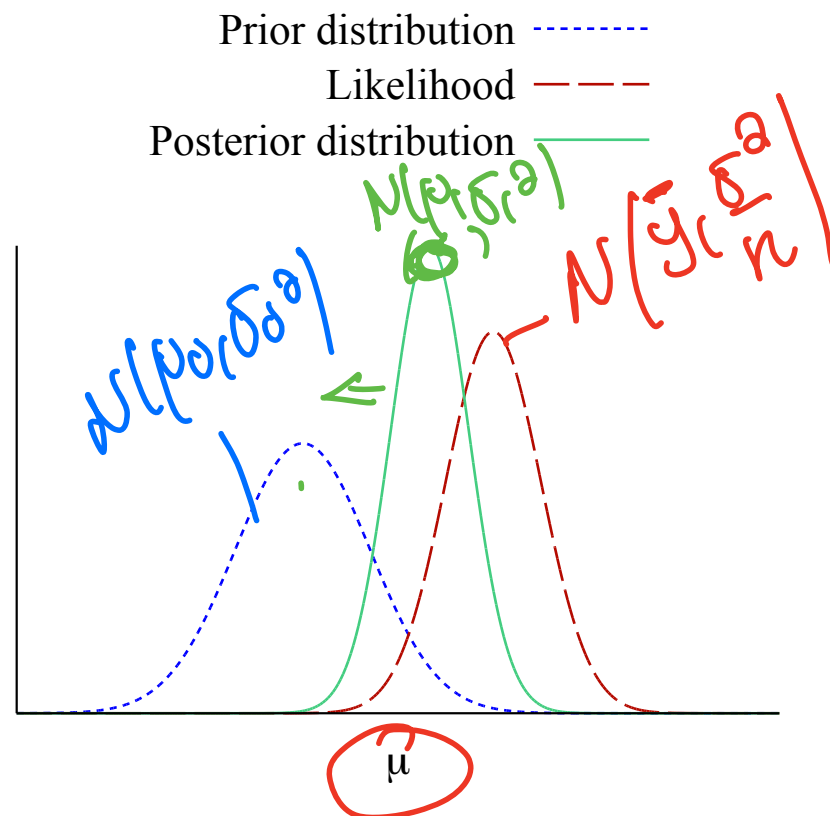
$$\mathcal{N}\left(\bar{y}, \left(\frac{\sigma^2}{n}\right)\right) \times \mathcal{N}\left(\mu_0, \sigma_0^2\right)$$

Informative and uninformative prior distributions

Normal example with known variance

- An informative prior distribution is strongly peaked around some value.
- Prior changes its value over the range of the likelihood.
- Posterior is shifted relative to likelihood.

$$\sigma_0^2 = 0.5, \frac{\sigma^2}{n} = 0.25$$



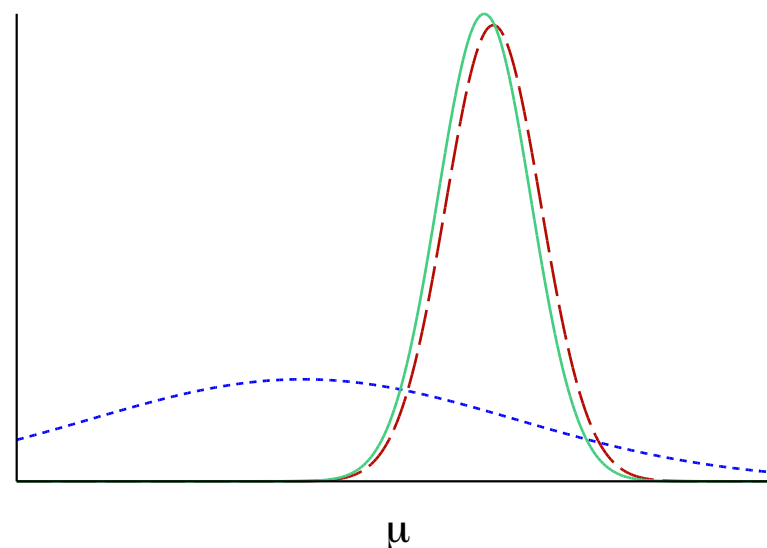
Informative and uninformative prior distributions

Normal example with known variance

- A weakly or slightly informative prior.
- Only changing gradually over the range of the likelihood.
- When the data provide a lot more information than the prior.
- Posterior is only slightly shifted relative to likelihood.

$$\sigma_0^2 = 5, \frac{\sigma^2}{n} = 0.25$$

Prior distribution -----
Likelihood - - - -
Posterior distribution ————



This prior is dominated by the likelihood and they give similar posterior.

Informative and uninformative prior distributions

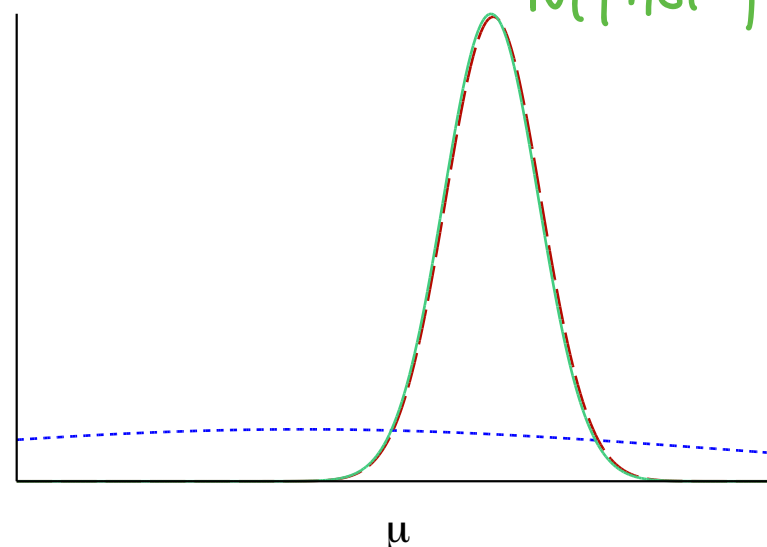
Normal example with known variance

- A very weakly informative prior, almost flat prior
- Almost flat over the range of the likelihood
- Posterior practically proportional to likelihood.

$$\sigma_0^2 = 20, \frac{\sigma^2}{n} = 0.25$$

$$p(\mu|y) \propto \text{likelihood} \propto N(\bar{y}, \frac{\sigma^2}{n})$$

Prior distribution
Likelihood --- $N(\bar{y}, \frac{\sigma^2}{n})$
Posterior distribution — $N(\mu, \sigma^2)$

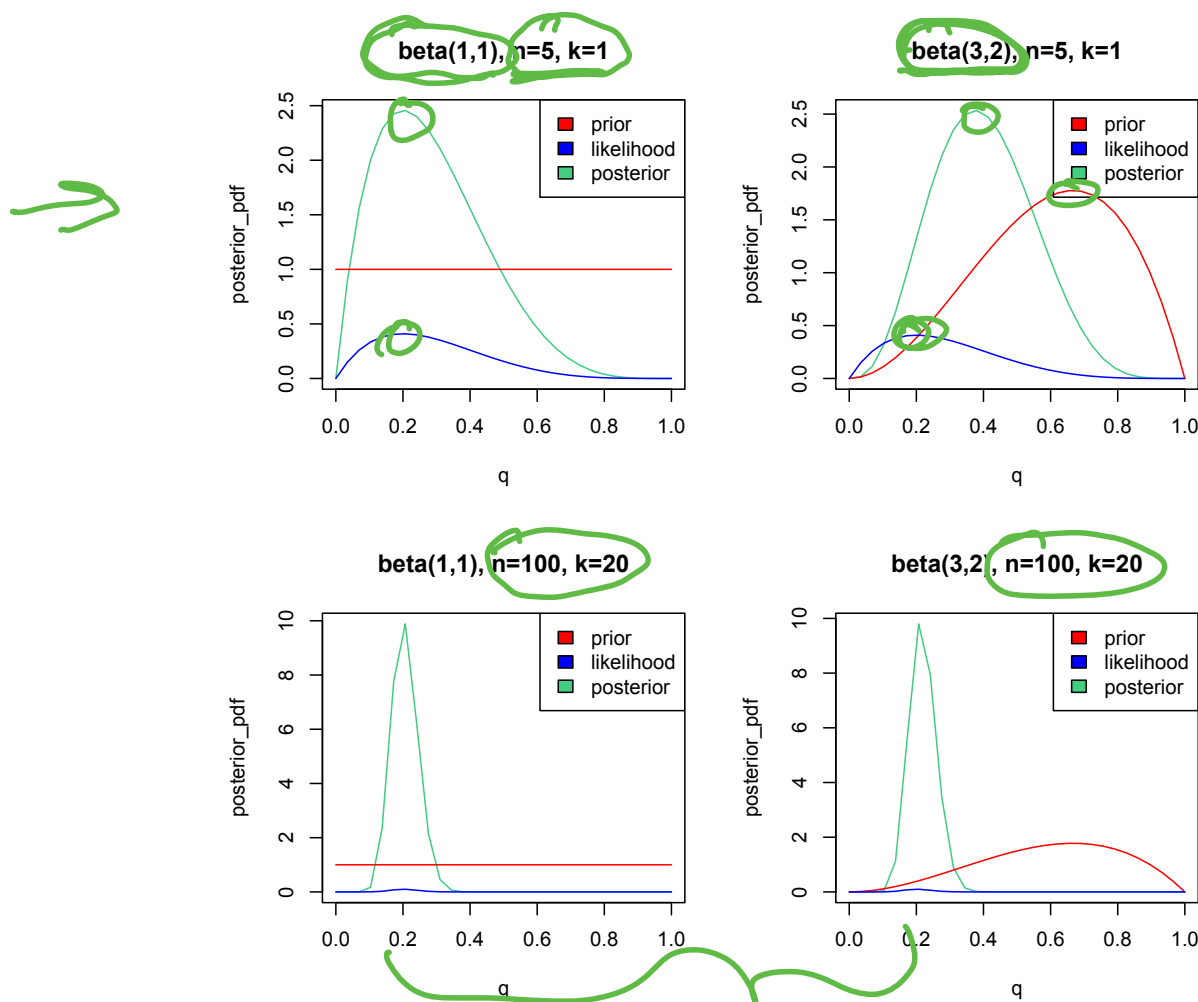


Beta prior/binomial data example

- Likelihood: $k \sim \text{binom}(n, q)$
- Prior on q : $p(q) \sim \text{beta}(\alpha, \beta)$, $q \in (0, 1)$.
- Posterior, $p(q|k) = \text{beta}(\alpha + k, \beta + n - k)$.
- When $\alpha = \beta = 1$, $q \sim \text{U}[0, 1]$ or $\text{beta}(1, 1)$.

Informative and uninformative prior distributions

Beta prior/binomial data example



more data lessens the dependence of the posterior on the prior.

Concept question: Normal/normal example

Question 2(a) from final exam Jan 2021

- 1 We have data $y = (y_1, \dots, y_n)$ from $N(\theta, \sigma^2)$, where $\sigma = 2$.
- 2 Prior distribution, $p(\theta) \sim N(0, \sigma_0^2)$.
- 3 **Question:** For an **uninformative** prior, do we need a large or small value for the prior standard deviation σ_0 ?

A larger standard deviation corresponds to a less informative prior or uninformative prior.

Concept question: Normal/normal example

Question 2(a) from final exam Jan 2023

- 1 Same normal/normal example with previous examples.
- 2 **Question:** As the prior distribution becomes less informative, what value does the posterior mean for θ approach? As the prior distribution becomes more informative, what value does the posterior mean for θ approach?

A larger value for σ_0 corresponds to a less informative prior. As σ_0 becomes larger, the posterior mean, μ_1 , approaches the MLE, \bar{y} .
As σ_0 becomes smaller, μ_1 approaches the prior mean.