

Bayesian inference

Lecture 14b (STAT 24400 F24)

1 / 16

Review: the Bayesian framework

$$\begin{cases} \theta \sim g(\cdot) & \leftarrow \text{PMF/density of prior distribution} \\ X_1, \dots, X_n \mid \theta \sim f(\cdot \mid \theta) \end{cases}$$

The *posterior distribution* is the conditional distribution of θ (conditioned on the observed data X_1, \dots, X_n).

Posterior PMF/density:

$$\begin{aligned} h(t \mid X_1, \dots, X_n) &= \frac{g(t)f(X_1, \dots, X_n \mid t)}{f(X_1, \dots, X_n)} \\ &= \left(\text{terms that don't depend on } t \right) \cdot g(t) \cdot \prod_{i=1}^n f(X_i \mid t) \end{aligned}$$

2 / 16

Example: Exponential with gamma prior

Suppose the data is drawn from an exponential distribution:

$$X_1, \dots, X_n \mid \lambda \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$$

Our prior on the parameter λ is:

$$\lambda \sim \text{Gamma}(k, r) \leftarrow \text{shape } k > 0, \text{ rate } r > 0$$

Some facts about $\text{Gamma}(a, b)$:

- mean = $\frac{a}{b}$, mode = $\frac{a-1}{b}$ (if $a > 1$)
- For integer k : $\text{Gamma}(k, b)$ is the distribution of a sum of k independent $\text{Exponential}(b)$ r.v.'s
- So by CLT, if a is large, $\text{Gamma}(a, b) \approx N\left(\frac{a}{b}, \frac{a}{b^2}\right)$

3 / 16

Example: Exponential (cont.)

Back to the example: $X_1, \dots, X_n \mid \lambda \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda)$

- Prior density $g(t) = \frac{r^k}{\Gamma(k)} t^{k-1} e^{-rt}$, $t \geq 0$
- Data density $f(x \mid t) = te^{-tx}$, $x \geq 0$
- Posterior density: $h(t \mid X_1, \dots, X_n) = \frac{g(t)f(X_1, \dots, X_n \mid t)}{f(X_1, \dots, X_n)}$

$$\begin{aligned} h(t \mid X_1, \dots, X_n) &= \left(\text{terms that don't depend on } t \right) \cdot \frac{r^k}{\Gamma(k)} t^{k-1} e^{-rt} \cdot \prod_{i=1}^n te^{-tX_i}, \quad t \geq 0 \\ &= \left(\text{terms that don't depend on } t \right) \cdot t^{k+n-1} e^{-(r+\sum_i X_i)t}, \quad t \geq 0 \end{aligned}$$

\rightsquigarrow the posterior distribution is $\text{Gamma}(k+n, r+\sum_i X_i)$

4 / 16

The Bayesian framework (point estimation via posterior)

The posterior gives a distribution of θ (given observed data).

What if we want a “point estimate”, i.e. a single value that is a good estimate for θ ?

Two standard options:

- Posterior mean:

$$\hat{\theta} = \mathbb{E}(\theta \mid X_1, \dots, X_n) \quad \leftarrow \mathbb{E}(\cdot) \text{ with respect to posterior } h(\cdot \mid X_1, \dots, X_n)$$

- Posterior mode (MAP):

$$\hat{\theta} = \operatorname{argmax}_{t \in \Theta} h(t \mid X_1, \dots, X_n)$$

5 / 16

Example: Exponential

Our model:

$$\begin{cases} \lambda \sim \text{Gamma}(k, r) \\ X_1, \dots, X_n \mid \lambda \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda) \end{cases}$$

Posterior:

$$\lambda \mid X_1, \dots, X_n \sim \text{Gamma}(k + n, r + \sum_{i=1}^n X_i)$$

Recall for $\text{Gamma}(a, b)$: mean = $\frac{a}{b}$, mode = $\frac{a-1}{b}$ (if $a > 1$)

$$\Rightarrow \text{Posterior mean} = \frac{k + n}{r + \sum_i X_i}, \quad \text{Posterior mode (MAP)} = \frac{k + n - 1}{r + \sum_i X_i}$$

6 / 16

Construction of credible intervals

A $(1 - \alpha)$ **credible interval** I (calculated as a function of X_1, \dots, X_n) contains $(1 - \alpha)$ posterior probability:

$$\mathbb{P}(\theta \in I \mid X_1, \dots, X_n) = 1 - \alpha$$

There are various ways to construct a credible interval.

Two common options:

- Equal tailed interval
- High posterior density interval

(For a symmetric & unimodal distribution, these options are equivalent)

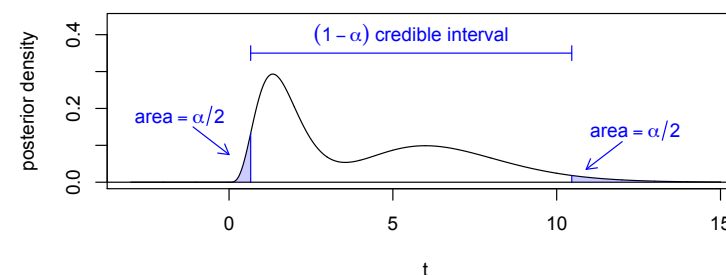
Remarks: the ideas used in the methods also apply to construction of frequentist confidence intervals with asymmetric and non-unimodal distributions.

7 / 16

Equal tailed credible intervals

- Equal tailed interval: our interval is

$$F_{\text{posterior}}^{-1}(\alpha/2) \leq \theta \leq F_{\text{posterior}}^{-1}(1 - \alpha/2).$$



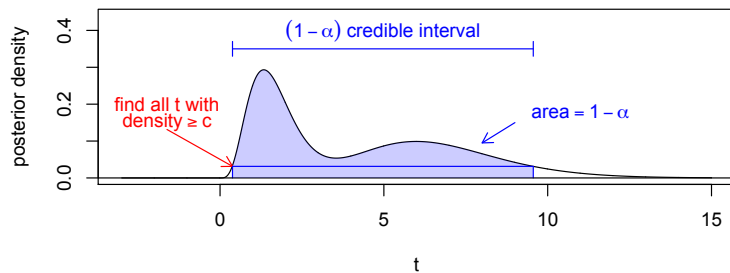
8 / 16

High posterior credible intervals

- High posterior density interval: our interval is given by

$$I = \{t : f_{\theta|X_1, \dots, X_n}(t | x_1, \dots, x_n) \geq c\}$$

where the density cutoff c is chosen so that $\text{prob.} = 1 - \alpha$



Note that this region I might not be a single interval!
(In the example above, if α is large, then I splits into two intervals)

9 / 16

Back to Example: Exponential (credible interval)

Our model:

$$\begin{cases} \lambda \sim \text{Gamma}(k, r) \\ X_1, \dots, X_n | \lambda \stackrel{\text{iid}}{\sim} \text{Exponential}(\lambda) \end{cases}$$

Posterior:

$$\lambda | X_1, \dots, X_n \sim \text{Gamma}(k + n, r + \sum_i X_i)$$

Equal-tailed credible interval:

$$F_{\text{Gamma}(k+n, r+\sum_i X_i)}^{-1}(\alpha/2) \leq \lambda \leq F_{\text{Gamma}(k+n, r+\sum_i X_i)}^{-1}(1 - \alpha/2)$$

write $F_{\text{Gamma}(a,b)}$ for the CDF of $\text{Gamma}(a, b)$

Note: the high posterior density interval is more complex to compute (omitted)

10 / 16

Example: exponential (normal approx. credible interval)

- Recall the fact about $\text{Gamma}(a, b)$ — For integer a :
 $\text{Gamma}(a, b)$ is the distribution of a sum of a indep. $\text{Exponential}(b)$ r.v.'s.
So by the CLT, for large integer a , $\text{Gamma}(a, b) \approx N(\frac{a}{b}, \frac{a}{b^2})$

- \Rightarrow for any $x \in (0, \infty)$, the CDF $F_{\text{Gamma}(a,b)}(x) \approx \Phi\left(\frac{x - a/b}{\sqrt{a/b}}\right)$
- \Rightarrow for any $t \in (0, 1)$,

$$F_{\text{Gamma}(a,b)}^{-1}(t) \approx F_{N(\frac{a}{b}, \frac{a}{b^2})}^{-1}(t) = \frac{a}{b} + \Phi^{-1}(t) \cdot \frac{\sqrt{a}}{b}$$

For example, for $t = 1 - \alpha/2$,

$$F_{\text{Gamma}(a,b)}^{-1}(1 - \alpha/2) \approx \frac{a}{b} + \Phi^{-1}(1 - \alpha/2) \cdot \frac{\sqrt{a}}{b} = \frac{a}{b} + z_{\alpha/2} \cdot \frac{\sqrt{a}}{b}$$

11 / 16

Example: exponential (Bayesian credible interval vs frequentist conf. interval)

Therefore, the $(1 - \alpha)$ (equal-tailed) credible interval

$$F_{\text{Gamma}(k+n, r+\sum_i X_i)}^{-1}(\alpha/2) \leq \lambda \leq F_{\text{Gamma}(k+n, r+\sum_i X_i)}^{-1}(1 - \alpha/2)$$

is approximately equal to:

$$\approx \frac{k + n}{r + \sum_i X_i} \pm z_{\alpha/2} \cdot \frac{\sqrt{k + n}}{r + \sum_i X_i}$$

If n is large (while k & r are constant), this credible int. is

$$\approx \frac{n}{\sum_i X_i} \pm z_{\alpha/2} \cdot \frac{\sqrt{n}}{\sum_i X_i} = \underbrace{\frac{1}{\bar{X}} \pm z_{\alpha/2} \cdot \frac{1}{\sqrt{n} \cdot \bar{X}}}}_{\text{= frequentist interval (using asymp. normality of the MLE)}}$$

12 / 16

Bayes risk

In the Bayesian framework, what's the best way to choose an estimator $\hat{\theta}$ to minimize squared loss?

At a *fixed* parameter value θ , the MSE is

$$\mathbb{E}((\hat{\theta} - \theta)^2) \leftarrow \mathbb{E}(\cdot) \text{ with respect to } X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(\cdot | \theta)$$

In a Bayesian framework, should also account for the random distrib. of θ :

$$\text{Bayes risk} = \mathbb{E}((\hat{\theta} - \theta)^2) \leftarrow \mathbb{E}(\cdot) \text{ with respect to } \begin{cases} \theta \sim g(\cdot) \\ X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} f(\cdot | \theta) \end{cases}$$

13 / 16

Bayes rule (for squared loss)

The **Bayes rule** is the estimator $\hat{\theta}$ (i.e., the function $\hat{\theta}(X_1, \dots, X_n)$) that minimizes Bayes risk

For squared loss:

$$\mathbb{E}((\hat{\theta} - \theta)^2) = \mathbb{E} \left(\mathbb{E}((\hat{\theta} - \theta)^2 | X_1, \dots, X_n) \right)$$

$\mathbb{E}(\cdot)$ with respect to marginal distrib. of X_1, \dots, X_n
 $\mathbb{E}(\cdot)$ with respect to posterior distrib. of $\theta | X_1, \dots, X_n$

This is minimized by $\hat{\theta} =$ posterior mean

14 / 16

Bayes rule (proof for squared loss)

Why does the posterior mean minimize $\mathbb{E}((\hat{\theta} - \theta)^2)$?

- For any random variable T and any constant t ,

$$\mathbb{E}((T - t)^2) = \text{Var}(T - t) + (\mathbb{E}(T - t))^2 = \text{Var}(T) + (\mathbb{E}(T) - t)^2$$

$\Rightarrow \mathbb{E}((T - t)^2)$ is minimized by choosing $t = \mathbb{E}(T)$

- For any random variables T and S , and any function $t(S)$,

$$\mathbb{E}((T - t(S))^2 | S) \text{ is minimized by choosing } t(S) = \mathbb{E}(T | S)$$

$\Rightarrow \hat{\theta} = \mathbb{E}(\theta | X_1, \dots, X_n) =$ posterior mean is the estimator that minimizes

$$\mathbb{E}((\hat{\theta} - \theta)^2 | X_1, \dots, X_n)$$

i.e., the expected squared error conditional on the data

$\Rightarrow \hat{\theta}$ must minimize $\mathbb{E}((\hat{\theta} - \theta)^2)$.

15 / 16

Other definitions of Bayes risk

We can generalize this to other loss functions $\text{loss}(\hat{\theta}, \theta)$, for example:

- Absolute loss:

$$\text{Bayes risk} = \mathbb{E}(|\hat{\theta} - \theta|)$$

\rightsquigarrow minimized by $\hat{\theta} =$ posterior median

- 0/1 loss: \leftarrow for the case of a discrete prior (& so the posterior is discrete)

$$\text{Bayes risk} = \mathbb{E}(\mathbb{1}_{\hat{\theta} \neq \theta}) = \mathbb{P}(\hat{\theta} \neq \theta)$$

\rightsquigarrow minimized by $\hat{\theta} =$ posterior mode

16 / 16