

STAT 8700 Homework 3

Brian Detweiler

Friday, September 16, 2016

1. Suppose we have a population described by a Normal Distribution with known variance $\sigma^2 = 1600$ and unknown mean μ . 4 observations are collected from the population and the corresponding values were: 940, 1040, 910, and 990.

```
y.bar <- mean(940, 1040, 910, 990)
y.bar
```

```
## [1] 940
```

(a) If we choose to use a Normal(1000, 200^2) prior for θ , and the posterior distribution for θ by hand.

We just need to find the likelihood $p(y|\sigma^2)$ for y_1, y_2, y_3, y_4 . Using the data, we have

$$\begin{aligned} p(y|\sigma^2) &= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(y_i - \mu)^2} \\ &= \prod_{i=1}^4 \frac{1}{\sqrt{\sigma^2}\sqrt{2\pi}} e^{-\frac{1}{2\cdot\sigma^2}(y_i - \mu)^2} \\ &= \prod_{i=1}^4 \frac{1}{\sqrt{\sigma^2}\sqrt{2\pi}} e^{-\frac{1}{2\cdot\sigma^2}(y_i - \mu)^2} \\ &= \left[\frac{1}{\sqrt{\sigma^2}\sqrt{2\pi}} \right]^4 e^{-\frac{1}{2\cdot\sigma^2} \sum_{i=1}^4 (y_i - \mu)^2} \\ &= (\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2} \end{aligned}$$

Now, let $v = \frac{1}{n} \sum_{i=1}^4 (y_i - \mu)^2$.

We now have

$$p(y|\sigma^2) \propto (\sigma^2)^{-\frac{n}{2}} e^{-\frac{n\cdot v}{2\sigma^2}}$$

This takes the shape of an Inverse-Gamma distribution.

$$\begin{aligned} Y &\sim \text{Inv-Gamma}(\alpha, \beta) \\ p(y|\alpha, \beta) &= \frac{\beta^\alpha}{\Gamma(\alpha)} y^{-(\alpha+1)} e^{-\frac{\beta}{y}} \end{aligned}$$

Now, let $\alpha = \frac{\nu}{2}, \beta = \frac{1}{2}$. This gives us

$$\begin{aligned} Y &\sim \text{Inv-Gamma}\left(\frac{\nu}{2}, \frac{1}{2}\right) \\ &= \frac{1}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})} y^{-(\frac{\nu}{2}+1)} e^{-\frac{1}{2y}}, y > 0 \end{aligned}$$

which is known as the Inverse-Chi-Square Distribution

■

(b) Find, by hand, a 95% credible interval for θ .

A 95% CI for θ is given by evaluating $p(y|\theta)$ at $y = 0.025$ and $y = 0.975$, with $\nu = 4$ degrees of freedom.

$$\begin{aligned} p(0.025; \theta) &= \frac{1}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})} y^{-\left(\frac{\nu}{2}+1\right)} e^{-\frac{1}{2y}}, y > 0 \\ &= \frac{1}{2^{\frac{4}{2}} \Gamma(\frac{4}{2})} (0.025)^{-\left(\frac{4}{2}+1\right)} e^{-\frac{1}{2(0.025)}} \\ &= \frac{1}{4} (0.025)^{-3} e^{-\frac{1}{0.05}} \\ &\approx 0.000032978457959 \\ p(0.975; \theta) &= \frac{1}{2^{\frac{\nu}{2}} \Gamma(\frac{\nu}{2})} y^{-\left(\frac{\nu}{2}+1\right)} e^{-\frac{1}{2y}}, y > 0 \\ &= \frac{1}{2^{\frac{4}{2}} \Gamma(\frac{4}{2})} (0.975)^{-\left(\frac{4}{2}+1\right)} e^{-\frac{1}{2(0.975)}} \\ &= \frac{1}{4} 1.07891232152 e^{-\frac{1}{1.95}} \\ &\approx 0.161514323478 \end{aligned}$$

This gives us a 95% Credible Interval of (0.000032978457959, 0.161514323478).

■

2. The `normnp` function in the `Bolstad` package computes the posterior for the mean with a Normal prior. The function requires 4 inputs (in order): a vector containing the data, the prior mean, the prior standard deviation, and the population standard deviation. Suppose we consider a Normal population with a variance of 16, and we collect 15 observations from this population with values: 26.8, 26.3, 28.3, 28.5, 26.3, 31.9, 28.5, 27.2, 20.9, 27.5, 28.0, 18.6, 22.3, 25.0, 31.5.

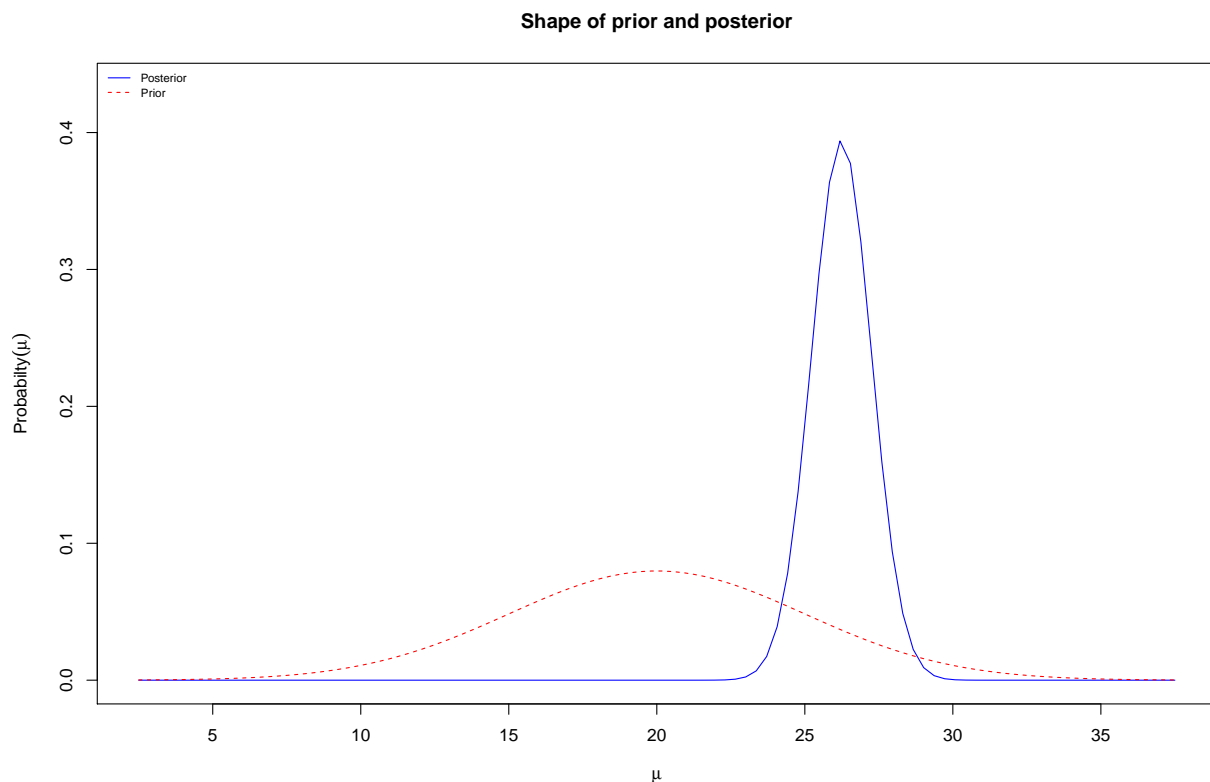
```
library(Bolstad)
var <- 16
obs <- c(26.8, 26.3, 28.3, 28.5, 26.3, 31.9, 28.5, 27.2, 20.9, 27.5, 28.0, 18.6, 22.3, 25.0, 31.5)
pop.st.dev <- sqrt(16)
```

(a) If we choose a $Normal(20, 25)$ prior, Use **R** to find the posterior distribution for the population mean.

```
prior.mu <- 20
prior.st.dev <- sqrt(25)

posterior <- normnp(obs, prior.mu, prior.st.dev, pop.st.dev)

## Known standard deviation :4
## Posterior mean           : 26.2404092
## Posterior std. deviation : 1.0114435
```



```
##
## Prob.    Quantile
## -----
## 0.005    23.6351035
## 0.010    23.8874398
## 0.025    24.2580164
## 0.050    24.5767327
## 0.500    26.2404092
## 0.950    27.9040857
## 0.975    28.2228020
## 0.990    28.5933786
## 0.995    28.8457149
```

(b) What are the posterior mean and variance?

The posterior mean is 26.2404092, and variance is 1.0230179.

(c) Find a 95% credible interval for the population mean.

A 95% credible interval for the population mean is found at the 0.025 and 0.975 quantiles, (24.2580164, 28.222802).

■

3. Suppose $y|\theta \sim \text{Poisson}(\theta)$, find the Jeffreys' prior density for θ . Find α and β for which the $\text{Gamma}(\alpha, \beta)$ density is a close match to the Jeffreys' prior.

Jeffrey's prior is given by $J(\theta) = \sqrt{I(\theta)}$, where $I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \ln p(y|\theta)\right]$.

The Poisson distribution we are interested in, is $p(y_n|\theta) = \theta^{\sum_{i=1}^n y_i} e^{-n\theta} \prod_{i=1}^n \frac{1}{y_i!}$.

So working through this by parts, we start with the natural log,

$$\begin{aligned} \ln \theta^{\sum_{i=1}^n y_i} e^{-n\theta} \prod_{i=1}^n \frac{1}{y_i!} &= \ln \frac{1}{y!} - \theta + y \ln \theta \\ &= \sum_{i=1}^n y_i \ln \theta - n\theta - \ln \sum_{i=1}^n y_i! \end{aligned}$$

Taking the first derivative with respect to θ , we get

$$\begin{aligned} \frac{\partial}{\partial \theta} \ln p(y_n|\theta) &= \frac{\partial}{\partial \theta} \sum_{i=1}^n y_i \ln \theta - n\theta - \ln \sum_{i=1}^n y_i! \\ &= \sum_{i=1}^n \frac{y_i}{\theta} - n - 0 \end{aligned}$$

Taking the second derivative with respect to θ , we get

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} \ln p(y_n|\theta) &= \frac{\partial}{\partial \theta} \sum_{i=1}^n \frac{y_i}{\theta} \\ &= - \sum_{i=1}^n y_i \frac{1}{\theta^2} \end{aligned}$$

Taking expectations,

$$\begin{aligned} -E\left[-\frac{y}{\theta^2} \middle| \theta\right] &= \frac{n\theta}{\theta^2} \\ &= \frac{n}{\theta} \end{aligned}$$

Finally, taking the square root to get the Jeffrey's prior, $J(I)$, we have

$$\begin{aligned} \sqrt{I(\theta)} &= \sqrt{\frac{n}{\theta}} \\ &\propto \sqrt{\frac{1}{\theta}} \\ &= \theta^{\frac{1}{2}} \end{aligned}$$

This comes closest to $\lim_{\beta \rightarrow 0} \text{Gamma}(\frac{1}{2}, \beta)$, though it is not a proper distribution.

■

4. Suppose we have multiple independent observations y_1, y_2, \dots, y_n from a $Poisson(\theta)$ distribution.

(a) Consider the conjugate Gamma prior. What values of the hyperparameters would lead to a flat (improper) prior distribution for θ ?

(b) Using a general $Gamma(\alpha, \beta)$ prior, derive the posterior distribution for θ . What is the required sufficient statistic needed from the data?

5. Derive the gamma posterior distribution (equation 2.15) for the Poisson model parameterized in terms of rate and exposure with conjugate prior distribution.

$$\begin{aligned}
 p(\theta|y) &\propto p(y|\theta)p(\theta) \\
 &\propto \left[\theta^{\sum_{i=1}^n y_i} e^{-(x_i)\theta} \right] \cdot \left[\frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta} \right] \\
 &\propto \left[\theta^{\sum_{i=1}^n y_i} e^{-(x_i)\theta} \right] \cdot \left[\theta^{\alpha-1} e^{-\beta\theta} \right] \\
 &= \theta^{(\alpha + \sum_{i=1}^n y_i - 1)} e^{-(\beta + \sum_{i=1}^n x_i)\theta}
 \end{aligned}$$

And thus we have the posterior as $\theta|y \sim Gamma(\alpha + \sum_{i=1}^n y_i, \beta + \sum_{i=1}^n x_i)$.

■

6. The table at the end of the assignment gives the number of fatal accidents and deaths on scheduled airline flights per year over a ten year period from 1976 to 1985.

(a) Assume that the number of fatal accidents in each year are independent with a $Poisson(\theta)$ distribution. Using a flat prior for θ , and the posterior distribution for θ based on the the 10 years of provided data. If you have a $Gamma(\alpha, \beta)$ distribution then the function `qgamma(q, shape=a, rate=b)` will return the q th quantile of the $Gamma(\alpha, \beta)$ distribution. Use this to find the ‘symmetric’ 95% credible interval for θ .

(b) Now assume that the number of fatal accidents in each year follow independent Poisson distributions with a constant rate and an exposure in each year proportional to the number of passenger miles flown. Again using a flat prior distribution for θ , determine the posterior distribution based on the data. (Estimate the number of passenger miles flown in each year by dividing the appropriate columns of table and ignoring round-off errors, death rate is per 100 million miles.) Give a 95% predictive interval for the number of fatal accidents in 1986 under the assumption that 8×10^{11} passenger miles are flown that year.

(c) Repeat (a) above, replacing ‘fatal accidents’ with ‘passenger deaths.’

(d) Repeat (b) above, replacing ‘fatal accidents’ with ‘passenger deaths.’

```
years <- c(1976:1985)
fatal.accidents <- c(24, 25, 31, 31, 22, 21, 26, 20, 16, 22)
passenger.deaths <- c(734, 516, 754, 877, 814, 362, 764, 809, 223, 1066)
death.rate <- c(0.19, 0.12, 0.15, 0.16, 0.14, 0.06, 0.13, 0.13, 0.03, 0.15)

airline.deaths <- as.data.frame(cbind(years, fatal.accidents, passenger.deaths, death.rate))
airline.deaths
```

##	years	fatal.accidents	passenger.deaths	death.rate
## 1	1976	24	734	0.19
## 2	1977	25	516	0.12
## 3	1978	31	754	0.15
## 4	1979	31	877	0.16
## 5	1980	22	814	0.14
## 6	1981	21	362	0.06
## 7	1982	26	764	0.13
## 8	1983	20	809	0.13
## 9	1984	16	223	0.03
## 10	1985	22	1066	0.15

