

HW_1

Brandon

9/1/2021

Problem 1

In this problem we will consider developing a Bayesian model for Poisson data; i.e., our observed data will consist of $Y_1, \dots, Y_n \stackrel{i.i.d}{\sim} \text{Poisson}(\lambda)$. Recall, a random variable Y is said to follow a Poisson distribution, with mean parameter λ if its pmf is given by

Note, the Poisson model is often used to analyze count data.

- For the Poisson model, identify the conjugate prior. This should be a general class of priors.
 - Here as Poisson is part of the exponential family, we can see that its conjugate family of priors are $\text{Gamma}(a, b)$
- Under the conjugate prior, derive the posterior distribution of $\lambda|y$. This should be a general expression based on the choice of the hyper-parameters specified in your prior.

b) Derive the posterior dist of $\lambda|y$.

consider $Y \sim \text{Poisson}(\lambda)$ and $\lambda \sim \text{Gamma}(a, b)$

then by defn we have

$$P(\lambda|y_1, \dots, y_n) = P(\lambda) \cdot \frac{P(y_1, \dots, y_n|\lambda)}{P(y_1, \dots, y_n)}$$

now

$$P(\lambda) = \frac{b^a}{\Gamma(a)} \lambda^{a-1} e^{-b\lambda} \quad \text{and} \quad P(y_1, \dots, y_n|\lambda) = \frac{e^{-n\lambda} \lambda^{\sum y_i}}{\lambda^n n!}$$

and

$$P(\lambda|y) \propto P(y|\lambda)P(\lambda)$$

plugging in we get

$$\begin{aligned} P(\lambda|y) &\propto \lambda^{a-1} e^{-b\lambda} e^{-n\lambda} \lambda^{\sum y_i} \cdot C(a, b, \lambda) \\ &\propto \underbrace{\lambda^{a-1+\sum y_i}}_a e^{-\lambda(b+n)} C(a, b, \lambda) \end{aligned}$$

group constants (non λ)

Gamma kernel

thus

$$\underline{P(\lambda|y) \propto \text{gamma}(a + \sum y_i, b + n)}$$

- c. Find the posterior mean and variance of $\lambda|y$. These should be general expressions based on the choice of the hyper-parameters specified in your prior.

c) Find the mean and variance of the posterior $\lambda|y$.

$$\text{as } P(\lambda|y) \sim \text{gamma}(a + \sum y_i, b+n)$$

then

$$E[\lambda|y] = \frac{a + \sum y_i}{b+n} = \frac{a}{b+n} + \frac{\sum y_i}{b+n}$$

$$= \frac{\frac{b}{b+n} \cdot \frac{a}{b}}{\frac{b}{b+n}} + \frac{n}{b+n} \frac{\sum y_i}{n}$$

here $\frac{a}{b}$ is the prior mean
and $\frac{1}{n} \sum y_i$ is the sample avg \bar{y}

$$\text{Var}(\lambda|y) = \frac{a}{b^2} = \frac{a + \sum y_i}{(b+n)^2}$$

- d. Obtain the MLE of λ . Develop and discuss a relationship that exists between the MLE and posterior mean identified in (c).

d) Obtain the MLE of λ , is there a relationship between $E[\lambda|y]$

$$p(y|\lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$$

$$p(\lambda|y) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} = \frac{1}{\prod_{i=1}^n y_i!} e^{-n\lambda} \lambda^{\sum y_i}$$

$$\ln(p(\lambda|y)) = -\ln(\prod y_i!) - n\lambda + \sum y_i \ln \lambda$$

$$\frac{d}{d\lambda} \ln(p) = -n + \frac{\sum y_i}{\lambda} \stackrel{\text{set}}{=} 0$$

$$\frac{\sum y_i}{\lambda} = n \rightarrow \lambda^{\text{MLE}} = \frac{1}{n} \sum y_i$$

$$\frac{d^2}{d\lambda^2} \ln(p) = -\frac{\sum y_i}{\lambda^2} < 0 \quad \text{confirms global MLE}$$

$$\lambda^{\text{MLE}} = \frac{\sum y_i}{n} \rightarrow E[\lambda|y] = \frac{a + \sum y_i}{b+n} = \frac{b}{b+n} \frac{a}{b} + \frac{n}{b+n} \frac{\sum y_i}{n}$$

we can see that the MLE is part of the Expectation of the posterior distribution.

- e. Write two separate R programs which can be used to both find a $(1 - \alpha)100\%$ equal-tailed credible interval and a $(1 - \alpha)100\%$ HPD credible interval for the Poisson model. These programs should take as arguments the following inputs: the observed data, prior hyper-parameters, and significance level.

```
# Takes a poisson model with a gamma prior and returns a (1-alpha) credible
# interval for the model
```

```
EQCI.poisson <- function(y , a , b, alpha){
  return(qgamma(c(alpha/2, 1-alpha/2),shape = a + sum(y),rate = b + NROW(y)))
}
```

```
# HPD Interval for a given h
HPD.poisson.h <- function(y ,h = .1, a = 1 , b = 1, plot = F, ...){
  apost <- a + sum(y)
  bpost <- b + NROW(y) #TO ensure we read vectors and dataframes the same
  if (apost >= 1) {
    mode <- (apost - 1)/(bpost)
    dmode <- dgamma(mode, shape = apost, rate = bpost)}
  else return("mode at 0: HPD not implemented yet")

  lt <- uniroot(f=function(x){
    dgamma(x,shape = apost, rate = bpost)/dmode - h},
    lower=0, upper=mode)$root

  ut <- uniroot(f=function(x){
    dgamma(x,shape = apost, rate=bpost)/dmode - h},
    lower=mode, upper= 100^100)$root

  coverage = pgamma(ut, shape = apost, rate = bpost) -
    pgamma(lt, shape = apost, rate = bpost)
  if (plot) {
    ld <- seq(0, 3, length=1000)
    plot(ld, dgamma(ld, shape = apost, rate = bpost),
      t="l",
      lty=1,xlab=expression(lambda),
      ylab="Posterior Density", ...)
    abline(h = h*dmode)

    segments(ut,0,ut,dgamma(ut,shape = apost,rate = bpost))
    segments(lt,0,lt,dgamma(lt,shape = apost,rate = bpost))

    title(bquote(paste("P(", .(round(lt, 2)), " < ", lambda, " < ",
      .(round(ut,2)), " | " , y, ") = ",
      .(round(coverage, 2)))))
  }

  return(c(lt,ut,coverage,h))
}

#Helper Function
Dev.HPD.poisson.h<-function(h, y, alpha){
  cov<-HPD.poisson.h(y, h, plot=F)[3]
  res<-(cov-(1-alpha))^2
  return(res)
}

# Returns HPD Interval for poisson prior and gamma posterior at a certain alpha
HPD_Interval.poisson <- function(y, a , b ,alpha, Plot = F, ...){
  h.final <- optimize(Dev.HPD.poisson.h, c(0,1), y = y, alpha = alpha)$minimum

  return(HPD.poisson.h(y, h.final, a, b , Plot))
}
```

```

}

# Generate test data
test = rpois(n=30, lambda=3)
# Interval
print(paste('Credible:',round(EQCI.poisson(test , 2, 1, .05)[1],4),
        '- ',round(EQCI.poisson(test , 2, 1, .05)[2],4)))

## [1] "Credible: 2.5665 - 3.8172"

print(paste('HPD:',round(HPD_Interval.poisson(test ,2, 1,.05)[1],4),
        '- ', round(HPD_Interval.poisson(test ,2, 1,.05)[2],4)))

## [1] "HPD: 2.5465 - 3.7943"

# Getting ranges to confirm that HPD is more restrictive
print(paste('Credible Range:',
        EQCI.poisson(test , 2, 1, .05)[2] - EQCI.poisson(test , 2, 1, .05)[1]))

## [1] "Credible Range: 1.25066570509891"

print(paste('HPD Range:',
        HPD_Interval.poisson(test ,2, 1,.05)[2] - HPD_Interval.poisson(test ,2, 1,.05)[1]))

## [1] "HPD Range: 1.24783286823883"

```

f. Find a data set which could be appropriately analyzed using the Poisson model. This data set should be of interest to you, and you should discuss, briefly, why the aforementioned model is appropriate; e.g., consider independence, identically distributed, etc. etc. You will also need to provide the source of the data.

- I chose the following data set on the number of births per woman in individual countries. This data is appropriate as each individual birth is independent of each other with similar opportunity. I found this data at the following link <https://data.worldbank.org/indicator/SP.DYN.TFRT.IN?end=2019&start=1960&view=chart>

g. Analyze the data set you have selected in (e). Provide posterior point estimates of λ , credible intervals, etc. etc. Your analysis should be accompanied by an appropriate discussion of your findings.

- Without any prior info we go ahead and set our prior to $\text{gamma}(a = 2, b = 1)$. Then our posterior distribution for $\lambda|x$ is given by: $\text{Gamma}(\text{shape} = 831, \text{rate} = 2)$. Using the MLE we see that λ can be estimated at 3.5. Using the posterior mean it is almost the same at 3.49. Using our HPD we can see that the posterior probability that $\lambda \in [3.2555, 3.7301]$ is 95%. For our data this means that for the year 1997, the probability that average number of births per woman is between $[3.2555, 3.7301]$ is 95%.

```

#Specifically looking at the year 1997
y97 <- y['1997']
# Calculating the posterior
apost <- 2 + sum(y97, na.rm = TRUE)
bpost <- 1 + NROW(y97)
print(paste("Posterior Distribution: Gamma(shape=",apost,'rate =',bpost,')'))

```

```
## [1] "Posterior Distribution: Gamma(shape= 831 rate = 238 )"
```

```
# We calculate the alpha = .05 credible and HPD intervals
print(paste('Credible Interval:',round(EQCI.poisson(y97 , 2, 1, .05)[1],4),
        '- ',round(EQCI.poisson(y97 , 2, 1, .05)[2],4)))
```

```
## [1] "Credible Interval: 3.2582 - 3.7329"
```

```
print(paste('HPD Interval:',round(HPD_Interval.poisson(y97 ,2, 1,.05)[1],4),
        '- ', round(HPD_Interval.poisson(y97 ,2, 1,.05)[2],4)))
```

```
## [1] "HPD Interval: 3.2555 - 3.7301"
```

```
# Calculating the MLE of Lambda and then the posterior mean
sprintf('MLE: %.2f',mean(y97$`1997`))
```

```
## [1] "MLE: 3.50"
```

```
sprintf('Posterior Mean: %.2f',(apost/bpost))
```

```
## [1] "Posterior Mean: 3.49"
```

Problem 2

An engineer takes a sample of 5 steel I beams from a batch and measures the amount (X) they sag under a standard load. The amounts in mm are 5.19, 4.72, 4.81, 4.87, 4.88. For this data set, it is known that the sag is $normal(\mu, \sigma^2)$, where the standard deviation $\sigma = .25$ is known. Use a $normal(0, 1)$ prior for μ .

- a. Find the posterior distribution of μ . Show all work!
 - The posterior distribution is given by $n(4.83358, \frac{1}{81})$

Problem 2

Sample of 5 steel beams $n=5$

observed data

X = Sag under a standard load

$$x = \{5.19, 4.72, 4.81, 4.87, 4.88\}$$

$$X \sim N(\mu, \frac{1}{16}) \text{ where } \tau = 0.25 = 1/4$$

$$\text{prior info } \mu \sim N(0, 1)$$

$$p(\mu) \sim N(0, 1)$$

a) Find the Posterior Distribution

$$p(x|\mu) \sim N(\mu, (\frac{1}{4})^2)$$

recall:

$$p(\mu|x) \propto p(x|\mu)p(\mu) \text{ and}$$

Subbing gives

$$p(\mu|x) \propto \left[\prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x_i - \mu)^2\right] \right] \left[\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}\mu^2\right] \right]$$

Combining constants and expanding

$$\propto \exp\left[-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2}\mu^2\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2 \right) - \frac{1}{2}\mu^2\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\sum_{i=1}^n x_i^2 - 2n\bar{x}\mu + n\mu^2 \right) - \frac{1}{2}\mu^2\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\sum_{i=1}^n x_i^2 - 2n\bar{x}\mu + (n+1)\mu^2 \right)\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\sum_{i=1}^n x_i^2 - 2n\bar{x}\mu + (n+1)\mu^2 \right)\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\sum_{i=1}^n x_i^2 - 2n\bar{x}\mu + (n+1)\mu^2 \right)\right]$$

$$\propto \exp\left[-\frac{1}{2} \left(\mu - \frac{n\bar{x}}{n+1} \right)^2 / \frac{1}{n+1} \right]$$

μ ————— $\frac{n\bar{x}}{n+1}$ ————— $\frac{1}{n+1}$
 normal kernel

$$p(\mu|x) \propto N\left(\frac{n\bar{x}}{n+1}, \frac{1}{n+1}\right)$$

From the given data $n=5$, $\bar{x}=4.894$ plugging in we get

$$p(\mu|x) \propto N\left(\frac{16(5)(4.894)}{16(5)+1}, \frac{1}{16(5)+1}\right)$$

$$\propto N(4.83358, 1/81)$$

b. Draw the density of the posterior and prior distribution of μ in the same figure.

```
# posterior and prior
par(mar = c(4, 4, 1, 1))
theta <- seq(-8, 8, by = 0.01)
plot(theta, dnorm(theta, mean = 4.83358, sd = sqrt(1/81)), t="l", lty=1,
      ylab="Posterior Density",
      xlab=expression(theta), col = 'dark red')
lines(theta, dnorm(theta, 0, 1), lty=2, col = 'dark blue')
legend(-8, 3.5, c("Posterior", "Prior"), lty=c(1, 2))
```

