

# PUBH 7440: Intro to Bayesian Analysis

## Homework from Week 3 — Due Feb 15th

**Stroke mortality in PA:** [Insert something about how stroke is the 4th leading cause of death and thus it's super important we study trends in stroke mortality...] Here, we want to look at county-level stroke mortality rates among those aged 65–74, 75–84, and 85+ from 2016 in the state of Pennsylvania. Specifically, our **data** consist of the number of deaths due to stroke,  $Y_{ia}$ , from county  $i$  and age-group  $a$  out of a population of size  $n_{ia}$ , where  $i = 1, \dots, N_s$  and  $a = 1, \dots, N_g$  ( $N_s = 67$  and  $N_g = 3$ ).

1. Assume  $Y_{ia} \sim \text{Pois}(n_{ia}\lambda_{ia})$  where  $\lambda_{ia} \sim \text{Gam}(Y_{0a}, n_{0a})$ . Using the pmf of the Poisson distribution and pdf of the gamma distribution below, show that this prior is a *(conditionally) conjugate prior* for  $\lambda_{ia}$  and write its full-conditional distribution. Note: Please use the parameterization of the gamma distribution below, *not* the one listed in Appendix A of CL3.

$$p(\lambda_{ia} | Y_{0a}, n_{0a}) = \frac{n_{0a}^{Y_{0a}}}{\Gamma(Y_{0a})} \lambda_{ia}^{Y_{0a}-1} \exp[-n_{0a}\lambda_{ia}]$$

2. I claim that we can interpret  $Y_{0a}$  and  $n_{0a}$  as the prior number of deaths and prior population size for age  $a$ , respectively. Aside from my infallibility as your professor, why does this interpretation make sense?
3. Some of our counties have small population sizes, thus we may want to consider the use of informative priors. To do this, we want to construct our priors such that they are consistent with (a) rate estimates that we might expect and (b) prior population sizes that respect the age distribution of Pennsylvania.

- Let  $\pi_a = \sum_i n_{ia} / \sum_a \sum_i n_{ia}$  denote the proportion of Pennsylvania's population belonging to each age group.
- Define  $\boldsymbol{\lambda}_0 = (\lambda_{01}, \lambda_{02}, \lambda_{03})$  to be the vector of our prior guesses at the age-specific mortality rates — 75, 250, and 1,000 deaths per 100,000 for ages 65–74, 75–84, and 85+, respectively (e.g.,  $\lambda_{01} = 75/100,000$ ). These are estimated from **data** from 2015 across the entire US.

If we want our priors to correspond to a county whose total 65+ population size is 10,000, explain why specifying  $n_{0a} = \pi_a \times 10,000$  and  $Y_{0a} = n_{0a}\lambda_{0a}$  could achieve this and the two above goals.

4. Due to data confidentiality issues, all counts  $Y_{ia} < 10$  have been suppressed in public-use data. How can we account for this in our analysis?
5. Write out the full hierarchical model; i.e., something like:

$$p([\text{all of the unknown parameters}] | [\text{all of the } \textit{known} \text{ data}]) \propto [\text{Likelihood}] \times [\text{All of the priors}] .$$

6. Using the data on Canvas and the code outline below, write a Gibbs sampler to fit this model.
7. After obtaining samples from the posterior distribution for each  $\lambda_{ia}$ , obtain samples from the posterior distribution of the county-specific age-adjusted mortality rates —  $\lambda_i = \sum_a \pi_a \lambda_{ia}$ . Using the last piece of code below, create a map of the posterior medians of the age-adjusted rates. Your estimates should be *similar to* (but not the same as) the estimates on CDC WONDER.
8. **OPTIONAL:** CDC WONDER's privacy protections have issues. For instance, we can modify our request to obtain the total number of deaths in the state of Pennsylvania for each age group. While we *could* use this information to improve our imputation step, what I want you to do is:
  - Keep track of the imputed values for each iteration of the Gibbs sampler.
  - Compare the posterior distribution for the total number of deaths (i.e., the true/uncensored  $Y_{ia}$ 's plus the imputed values) to the true total death counts. Is our approach overestimating the death counts?

```
#https://wonder.cdc.gov/controller/saved/D140/D34F844
rm(list=ls())
#First we read in the data and define a few things...
stroke=read.table('2016_PA_stroke_total.txt',sep='\t',
                   stringsAsFactors=FALSE,header=TRUE)
Ng=3          #three age groups
alabs=unique(stroke$Age.Group.Code)
Ns=67         #67 counties
clabs=unique(stroke$County)

#Next we organize things a bit...
Y=array(stroke$Deaths,dim=c(Ng,Ns))
n=array(stroke$Population,dim=c(Ng,Ns))

#####
#####
#per part 4, all Y's below 10 have been suppressed
#####
thres=10      #Suppression threshold
dY=!is.na(Y)  #0 for suppressed, 1 for observed
nsupp=apply(!dY,1,sum) #how many suppressed per age
#note: we do not know the true Y's
#      CDC did the suppression, not me

#####
#####
#insert your prior info here
```

```

#####
lam0=c(75,250,1000)/100000
n0=      #####
Y0=      #####
#####

#####
#####  

#initialize your Gibbs sampler here
nsims=10000
lami=array(dim=c(Ng,Ns,nsims))
for(a in 1:Ng){
  lami[a,,1]=  #####
  Y[a,!dY[a,]]=  #####
  #Note: the preceding line assumes we don't care
  #      what the posterior dist of the missing
  #      Y's looks like -- I'm just plugging the
  #      current guesses directly into my data vector
}

for(it in 2:nsims){
  for(a in 1:Ng){
    #####
    #####  

    #ADDRESS SUPPRESSED Y HERE
    #####  

    #####
    #####
    #####
    #####
    #####
    #ESTIMATE LAMBDA_{ia} HERE
    #####
    #####
    #####
  }
}

#####
#####
#Get posterior samples
#of the age-adjusted rates
#####
aalami=array(dim=c(Ns,nsims))
for(i in 1:Ns){
  aalami[i,]=  #####
}

```

```

#####
#####
#calculate the posterior medians
# of the age-adjusted rates
aa.med= #####
#####

#####
#####
#THE BELOW CODE SHOULD BE LEFT AS-IS!
#IT ASSUMES YOU NAMED
#THE POSTERIOR MEDIAN OF THE AGE-ADJUSTED RATES
#"aamed" USING THE CODE ABOVE,
#AND WILL CREATE A MAP "PAmmap.png"
#THAT WILL BE SAVED TO YOUR CURRENT DIRECTORY
#####

load('penn.rdata')
install.packages(c('maptools','RColorBrewer'))
library(maptools)
library(RColorBrewer)
ncols=7
cols=brewer.pal(ncols,'RdYlBu')[ncols:1]
tcuts=quantile(aa.med*100000,1:(ncols-1)/ncols)
tcolb=array(rep(aa.med*100000,each=ncols-1) > tcuts,
            dim=c(ncols-1,Ns))
tcol =apply(tcolb,2,sum)+1

png('PAmmap.png',height=520,width=1000)
par(mar=c(0,0,0,10),cex=1)
plot(penn,col=cols[tcol],border='lightgray',lwd=.5)
legend('right',inset=c(-.15,0),xpd=TRUE,
       legend=c(paste(
           c('Below',round(tcuts[-(ncols-1)],0),'Over'),
           c(' ',rep(' - ',ncols-2),' '),
           c(round(tcuts,0),round(tcuts[ncols-1],0)),sep=''))),
       fill=cols,title='Deaths per 100,000',bty='n',cex=1.5,
       border='lightgray')
dev.off()

```

## Problem 1

- $Y_{i2} \sim \text{Pois}(u_{i2} \lambda_{i2})$ ,  $\lambda_{i2} \sim \text{Gamma}(Y_{02}, u_{02})$

$i \rightarrow$  country number

$\ell \rightarrow$  age group

- $Y_{i2}$  = deaths due to deaths stroke

$u_{i2}$  = population

$\lambda_{i2}$  = death rate

- $P(Y_{i2} | u_{i2} \lambda_{i2}) = \frac{e^{-(u_{i2} \lambda_{i2})} (u_{i2} \lambda_{i2})^{Y_{i2}}}{(Y_{i2})!}$

- $P(\lambda_{i2} | Y_{02}, u_{02}) = \frac{u_{02}^{Y_{02}}}{\Gamma(Y_{02})} \cdot \lambda_{i2}^{Y_{02}-1} e^{-u_{02} \lambda_{i2}}$

Posterior:

$$P(\underline{\lambda_{i2}} | \underline{Y_{i2}}) \propto \frac{e^{-\underline{(u_{i2} \lambda_{i2})}} (\underline{u_{i2} \lambda_{i2}})^{\underline{Y_{i2}}}}{(\underline{Y_{i2}})!} \times$$

$$x \frac{\frac{\gamma_{02}}{n_{02}}}{\Gamma(\gamma_{02})} \cdot \gamma_{iz}^{\gamma_{02}-1} e^{-\frac{n_{02}\gamma_{iz}}{\gamma_{02}}} \propto$$

$$\frac{e^{-(\gamma_{iz})}}{\gamma_{iz}} \gamma_{iz}^{\gamma_{iz}} \propto \frac{\gamma_{02}^{-1} e^{(-n_{02}\gamma_{02})}}{\gamma_{iz}^{(\gamma_{02} + \gamma_{iz})-1} e^{-(n_{02} + n_{iz})\gamma_{iz}}} \propto (\gamma_{iz})!$$

This resembles a kernel of a gamma distribution, so,

we conclude that a posterior distribution of  $\gamma_{iz}$  is given by

$$\gamma_{iz} | \gamma_{iz} \sim \text{Gamma}(\gamma_{02} + \gamma_{iz}, n_{02} + n_{iz})$$

So, a full conditional distribution can be written as

$$p(\gamma_{iz} | \gamma_{iz}, \gamma_{02}, n_{02}, n_{iz}) =$$

$$= \frac{(\kappa_{02} + \kappa_{i2})^{\gamma_{i2} + \gamma_{02}}}{\Gamma(\gamma_{i2} + \gamma_{02})} \times \frac{(\gamma_{02} + \gamma_{i2} - 1)}{\gamma_{i2}} e^{-(\kappa_{02} + \kappa_{i2})x_{i2}}$$

## Problem 6

### Gibbs Sampler

- 1) want to estimate  $\lambda_{id}$  distribution.
  - have some observed  $Y_{id}$   
 $u_{id}$
  - some  $Y_{id}$  need to be imputed
- 2) we know that  $\lambda_{id} | Y_{id} \sim \text{Gamma}(Y_0 + Y_{id}, n_0 + n_{id})$ .

Gibbs sampler should be:  
 $n_0 + n_{id}$ .

  - 1) we have crude greekes of  
75, 200, 1000 death per 100,000
  - 2) so,  $\lambda_{i1} = \frac{75}{100,000}$  for each age group &  
for all countries  $i=1, 2, \dots, 67$   
just to initialize the value
  - 3) we need to update  $\lambda_{id}$
  - 4) But!!! to update  $\lambda_{id}$  we need a  
value of  $Y_{id}$ .

But!!! Some of  $Y_{id}$  are missing.
  - 5) when they are missing, they are

below 10

- b) we can sample  $Y_{iz}$  from  
 $Y_{iz} \sim \text{Pois}(n_{iz} \times \pi_{iz})$ , and truncate  
at 10.

? So, Gibbs sampler:

- Sample  $Y_{iz}$ , impute  $\alpha_{iz}$ , use  $\pi_{iz}$
- update  $n_{ot+hiz}$ ,  $Y_{oz} + Y_{iz}$
- use  $\rightarrow$  as parameters of Gamma,  
sample  $\pi_{iz}$ .