# 732A91: Lab 1
# Bayesian Learning

Sarah Alsaadi, Carles Sans Fuentes

April 13, 2017

---

## Question 1

Let $y_1, \ldots, y_n | \theta \sim \text{Bern}(\theta)$, and assume that you have obtained a sample with $s = 14$ successes in $n = 20$ trials. Assume a $\text{Beta}(\alpha_0, \beta_0)$ prior for $\theta$ and let $\alpha_0 = \beta_0 = 2$.

- Draw random numbers from posterior $\theta | y_1, \ldots, y_n \sim \text{Bern}(\alpha_0 + s, \beta_0 + f)$ and verify graphically that the posterior mean and standard deviation converges to the true mean $E[\theta] = \frac{\alpha_0 + s}{\alpha_0 + s + \beta_0 + f} \approx 0.66$ and true standard deviation $Var(\theta) = \frac{(\alpha_0+s)*(\beta_0+f)}{(\alpha_0+s+\beta_0+f)^2*(\alpha_0+s+\beta_0+f+1)} \approx 0.09$. Figure 1 and 2 shows how the mean and standard deviation of samples randomly drawn converges to the true mean and true standard deviation with larger samples.
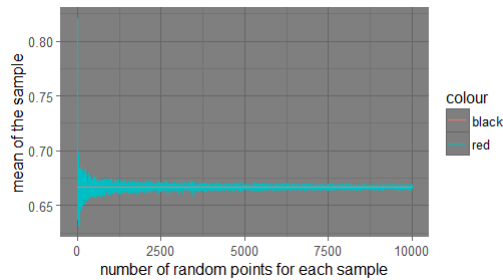


Figure 1: A plot representing the mean of 10000 samples where the sample size goes from 1 to 10000. The theoretical mean is plotted as a line.
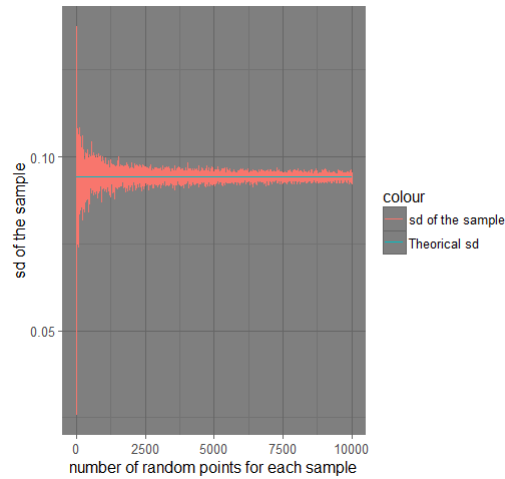
Figure 2: A plot representing the standard deviation of 10000 where the sample size goes from 1 to 10000. The theoretical standard deviation (sd) is plotted as a line.

- Use simulation (nDraws= 10000) to compute the posterior probability $Pr(\theta < 0.4|y)$ and compare with the exact value $Pr(\theta < 0.4|y) = 0.00397$. The simulated value is given below, we get 0.0036 which is close enough to the exact value.

```
1 pbeta(0.4,16,8)
2 [1] 0.003972563
3 simulated pbeta
4 [1] 0.0036
```

- Compute the posterior distribution of the log-odds $\phi = log(\frac{\theta}{1-\theta})$ by simulation with nDraws= 10000. Figure 3 below shows the histogram together with the kernel density of the data simulated from the posterior distribution of the log-odds $\phi = log(\frac{\theta}{1-\theta})$ with nDraws= 10000.
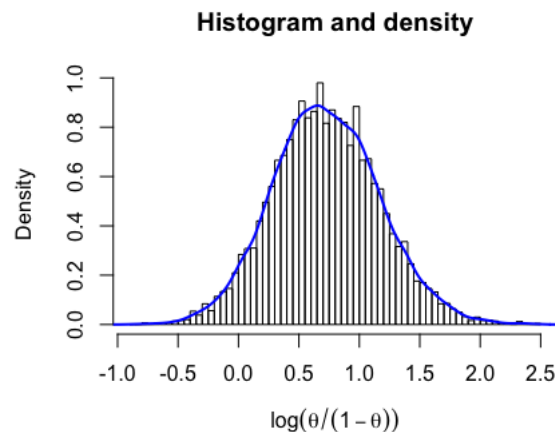


Figure 3: Histogram and kernel density of of the data simulated from the posterior distribution of the log-odds $\phi = log(\frac{\theta}{1-\theta})$ with nDraws= 10000.

## Question 2

Assume that you have asked 10 randomly selected persons about their monthly income (in thousands Swedish Krona) and obtained the following 10 observations
$(y_1, \ldots, y_{10}) = (14, 25, 45, 25, 30, 33, 19, 50, 34, 67)$. Assume $y_1, \ldots, y_n|\mu, \sigma^2 \sim logNormal(\mu, \sigma^2)$,

$\mu = 3.5$ and a non-informative prior $p(\sigma^2) \propto \frac{1}{\sigma^2}$. It can be shown that the posterior for $\sigma^2$ is $Inv - \chi^2(n, \tau^2)$ (scaled), where

$$\tau^2 = \frac{\sum_{i=1}^{n}(log(y_i) - \mu)}{n}$$

.

- Simulate 10000 draws from the posterior of $\sigma^2$ and compare with the theoretical $Inv - \chi^2(n, \tau^2) = Inv - \chi^2(10, 0.198)$. Figure 4 shows the histogram of the data simulated from $Inv - \chi^2(10, 0.198)$ and the density of the theoretical posterior distribution $Inv - \chi^2(10, 0.198)$. Since the simulated sample is so big, the histogram and the density is quite alike.
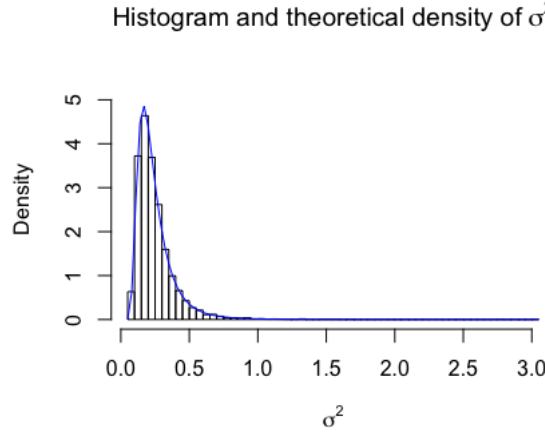
### Histogram and theoretical density of $\sigma^2$



Figure 4: Histogram of simulated data and theoretical density of the $\sigma^2$ with draws= 10000.

- The most common measure of income inequality is the Gini coefficient, G, where $0 \leq G \leq 1$. G= 0 means a completely equal income distribution, whereas G= 1 means complete income inequality. It can be shown that $G = 2\Phi(\frac{\sigma}{\sqrt{2}}) - 1$ when incomes follow a $logNormal(\mu, \sigma^2)$ distribution. Use the draws in a) to compute the posterior of the Gini coefficient for the current data set. Figure 5 shows the histogram of G and the kernel density. It shows that G is far from 1 so the income distribution seems to be closer to equal than inequal.

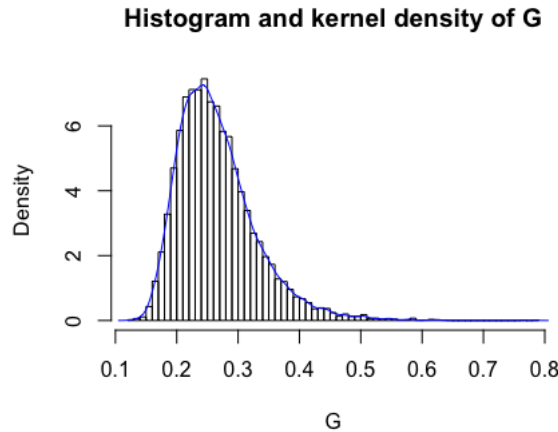### Histogram and kernel density of G



Figure 5: Histogram of simulated data and kernel density of G with draws= 10000.

- Use the posterior draws from b) to compute a 95% equal tail credible set for G. Also, do a kernel density estimate of the posterior of G and use it to compute a 95% Highest Posterior Density set for G. The Highest Posterior Density set (HPD) is a credible set that consists of

$\theta$-values having the Highest Posterior Density. Below we have a 95% equal tail credible set for G and and HPD. For a symmetric distribution, the intervals would be the same but since the distribution of G is slightly skewd, the intervals differ.

```
1 quantile(G,  probs = c(0.025, 0.975))
2 #    2.5%     97.5%
3 # 0.1739323 0.4152031
4 hdi(density(G),credMass=0.95)
5 #lower      upper
6 #0.1601422 0.3918670
```

# Question 3

The following data is the observed wind directions at a given location on 10 different days. The data recorded in radians:

$$(y_1, \ldots, y_{10}) = (-2.44, 2.14, 2.54, 1.83, 2.02, 2.33, -2.79, 2.23, 2.07, 2.02)$$

where $-\pi \leq y \leq \pi$.
Assume that the observations are independent and follow the von Mises distribution

$$p(y|\mu, \kappa) = \frac{exp(\kappa * cos(y - \mu))}{2\pi I_0(\kappa)}$$

$-\pi \leq y \leq \pi$.
Furthermore, assume that $\mu = 2.39$ and let $\kappa \sim \text{Exp}(\lambda = 1)$ apriori.

- Plot the posterior distribution of $\kappa$ for the data given above. The posterior is given by the following:

$$p(\kappa|\mu, y) = exp(-\kappa) \prod_{i=1}^{10} \frac{exp(\kappa * cos(y_i - 2.39))}{2\pi I_0(\kappa)}$$

Figure 6 shows the posterior of $\kappa$.
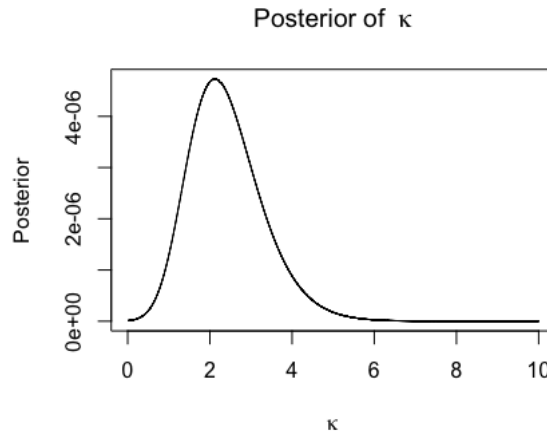


Figure 6: Posterior of $\kappa$ .

- The approximate point where the posterior has its mode is $(2.125, 4.7 * 10^{-6})$.

# Contributions

All results and comments presented have been developed and discussed together by the members of the group.

# Appendix

## Question 1

```
1
2  install.packages("geoR")
3  install.packages("HDInterval")
4  library("geoR")
5  library("HDInterval")
6  #lab 1 Bayesian learning
7
8  #1a
9  x<-rbeta(10000,2+14,2+6)
10 h<-hist(x,breaks = seq(0,1,0.02))
11 xfit<-seq(min(x),max(x),length=100)
12 yfit<-dbeta(xfit,16,8)
13 yfit <- yfit*diff(h$mids[1:2])*length(x)
14 lines(xfit, yfit, col="blue", lwd=2)
15
16 #1b
17 pbeta(0.4,16,8)
18 y<-x[x<0.4]
19 prob<-length(y)/length(x)
20
21
22 #1c
23 xodds<-log(x/(1-x))
24 hist(xodds,breaks = 50,prob=TRUE, main = 'Histogram and density', xlab = expression(paste(log(
       theta/(1-theta)))))
25 lines(density(xodds), col="blue", lwd=2)
26
27 #2a
28
29 data<-c(14,25,45,25,30,33,19,50,34,67)
30
31 tao2<-function(data)
32 {
33   sum((log(data)-3.5)^2)/length(data)
34
35 }
36
37 tao<-tao2(data=data)
38
39 sigma2<-rinvchisq(10000,df=10,scale=tao)
40 hist(sigma2,breaks = 100, prob=TRUE)
41 x<-seq(from=0, to=10000, by=0.001)
42 hx<-dinvchisq(x,df=10,scale=tao)
43
44 sigma2.histogram = hist(sigma2, breaks = 100, freq = F)
45 sigma2.ylim.normal = range(0, sigma2.histogram$density, dinvchisq(sigma2,df=10,scale=tao), na.
       rm = T)
46 hist(sigma2, breaks = 100, freq = F, ylim = c(0, 5.5), main=expression("Histogram and
       theoretical density of" ~ sigma^2), xlab = expression(paste(sigma^2)), ylab = 'Density')
47 curve(dinvchisq(x,df=10,scale=tao), add = T,col="blue")
48
49 #2.b
50 sigma<-sqrt(sigma2)/sqrt(2)
51 G<-2*pnorm(sigma,0,1)-1
52 hist(G,freq=F,breaks=50, main='Histogram and kernel density of G')
53 lines(density(G),col="blue")
54
55 #2.c
56 quantile(G,  probs = c(0.025, 0.975))
57 #    2.5%     97.5%
58 # 0.1739323 0.4152031
59
60 hdi(density(G),credMass=0.95)
61
62 #lower     upper
63 #0.1601422 0.3918670
64
65 #3.a
66
67 windradians<-c(-2.44,2.14,2.54,1.83,2.02,2.33,-2.79,2.23,2.07,2.02)
68 posterior<-function(k)
69 {
70   return(exp(-k)*prod(exp(k*cos(windradians-2.39))/(2*pi*besselI(x=k,nu=0))))
71
72 }
73
74
75 values<-sapply(k,posterior)
76 plot(k,values,type="l",main="Posterior of " ~ kappa, xlab=expression(paste(kappa)), ylab="
       Posterior")
77
```

```
78 #3.b
79 max(values)
80 #4.727694e-06
81 k[which.max(value)]
82 #2.125
```