

Computer Lab 1

732A54 - Bayesian Learning

Jooyoung Lee - joole336

Zuxiang Li - zuxli371

April 10, 2020

1. Bernoulli ... again.

Let $y_1, \dots, y_n | \theta \text{ Bern}(\theta)$, and assume that you have obtained a sample with $s = 5$ successes in $n = 20$ trials. Assume a $\text{Beta}(\alpha_0, \beta_0)$ prior for θ and let $\alpha_0 = \beta_0 = 2$.

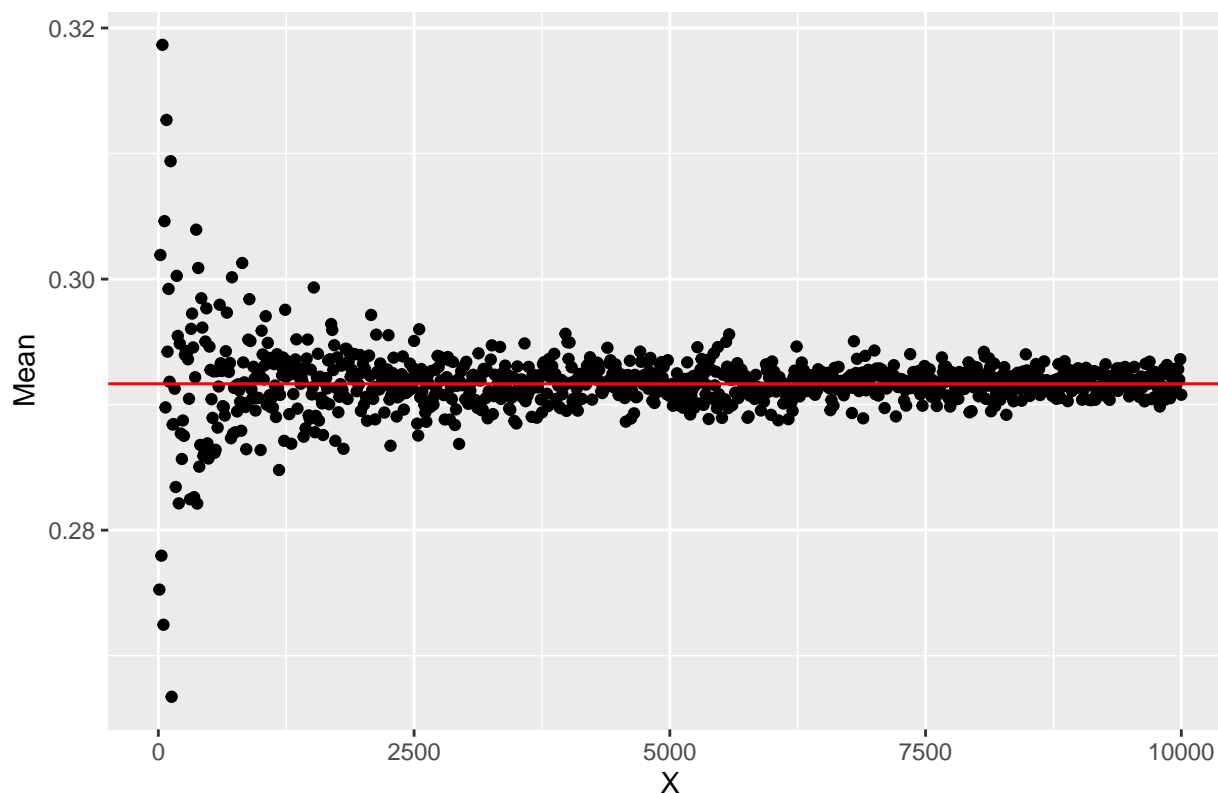
- (a) Draw random numbers from the posterior $\theta | y \text{ Beta}(\alpha_0 + s, \beta_0 + f)$, $y = (y_1, \dots, y_n)$, and verify graphically that the posterior mean and standard deviation converges to the true values as the number of random draws grows large.

```
## true mean= 0.2916667
```

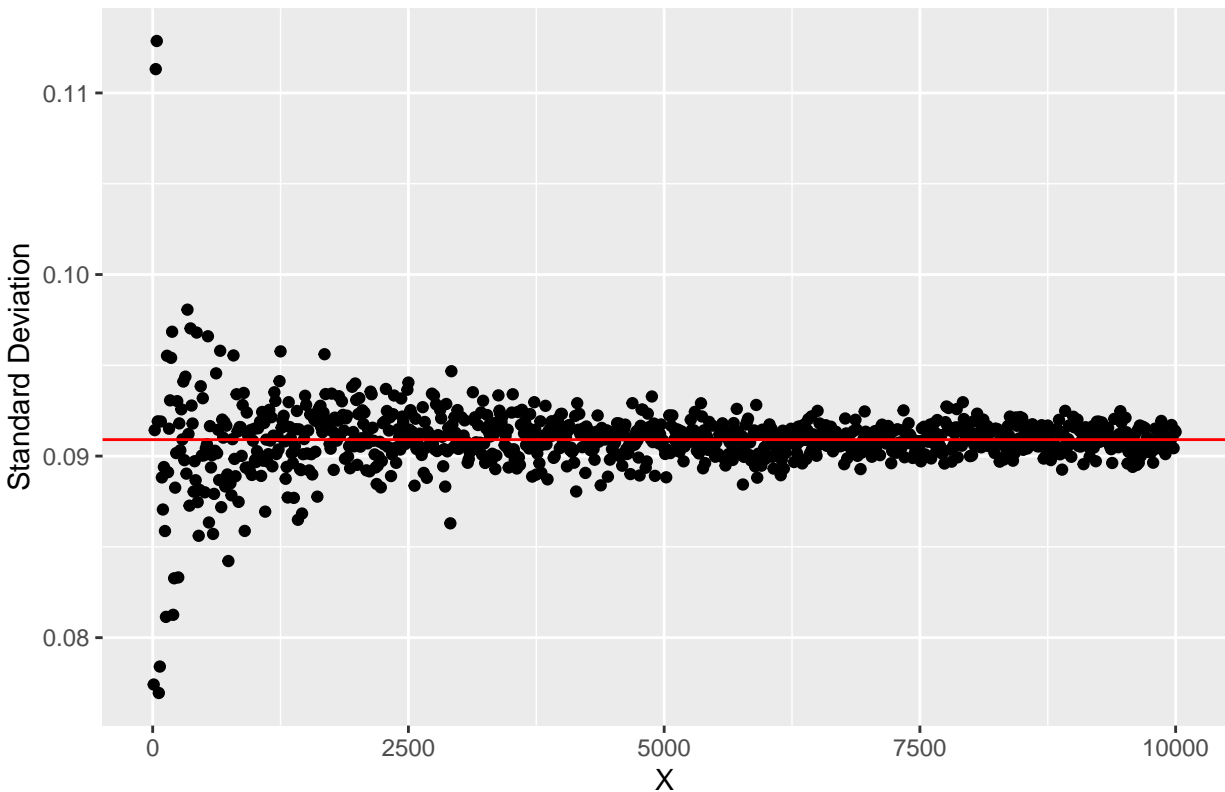
```
##
```

```
## true standard deviation= 0.09090593
```

Posterior mean with increasing number of random draws



Posterior standard deviation with increasing number of random draws



As we can see from the plots above, for both mean and stanard deviation, at first the values are distributed around two side of the red line. But as the number of random draws grows, the posterior mean and stanard deviation converges to the red line, which is the true mean and stanard deviation in this case.

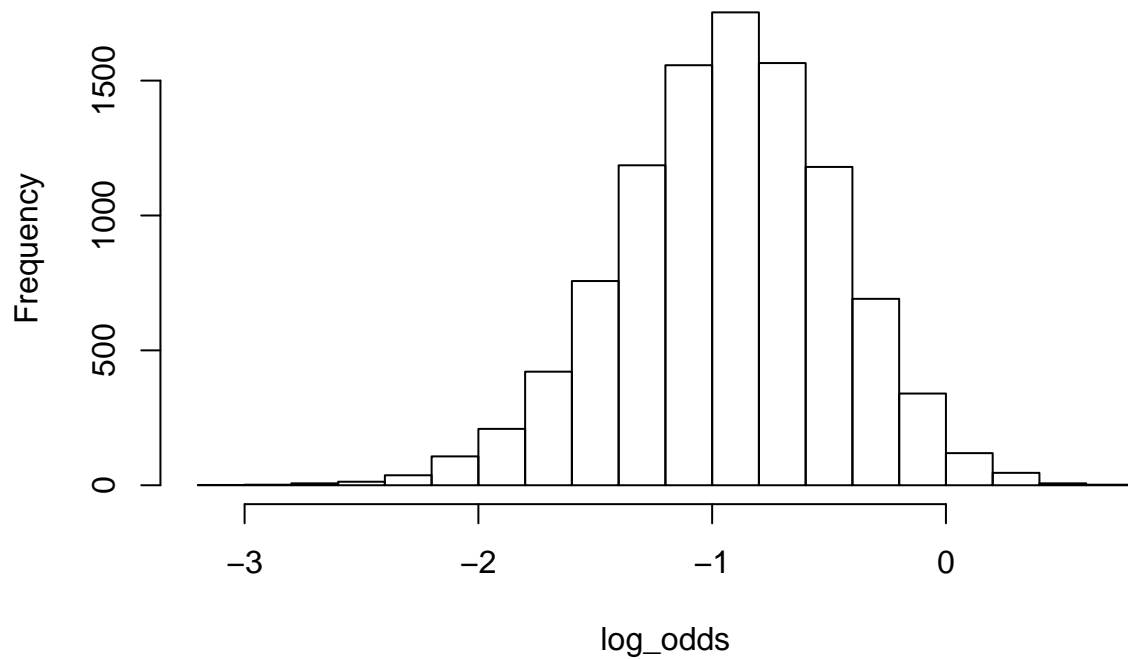
- (b) Use simulation ($n\text{Draws} = 10000$) to compute the posterior probability $Pr(\theta > 0.3|y)$ and compare with the exact value [Hint: `pbeta()`].

```
## sample prob= 0.4402
##
## true prob= 0.4399472
```

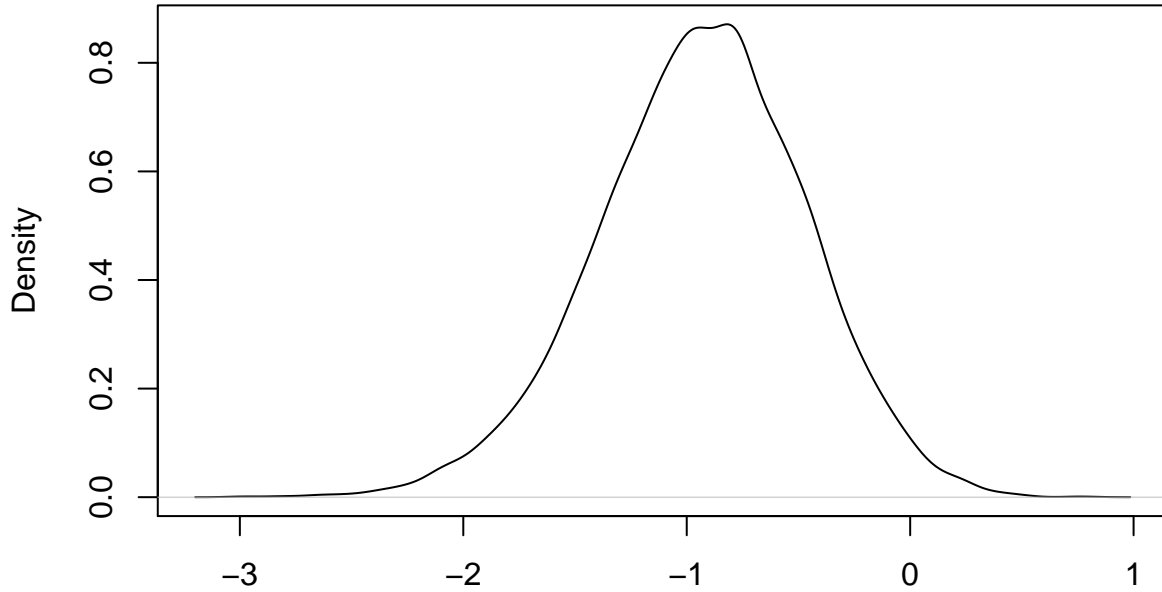
We use the α and β from (a) with function `rbeta()`, generated 10000 draws. Then get the number of thetas which are greater than 0.3 and divided by the total number. Using `pbeta()` to get the true probability.

- (c) Compute the posterior distribution of the log-odds $\phi = \log \frac{\theta}{1-\theta}$ by simulation ($n\text{Draws} = 10000$). [Hint: `hist()` and `density()` might come in handy].

Histogram of log_odds



density.default(x = log_odds)



N = 10000 Bandwidth = 0.06596

2. Log-normal distribution and the Gini coefficient.

Assume that you have asked 10 randomly selected persons about their monthly income (in thousands Swedish Krona) and obtained the following ten observations: 44, 25, 45, 52, 30, 63, 19, 50, 34 and 67. A common model for non-negative continuous variables is the log-normal distribution. The log-normal distribution

$\log N(\mu, \sigma^2)$ has density function

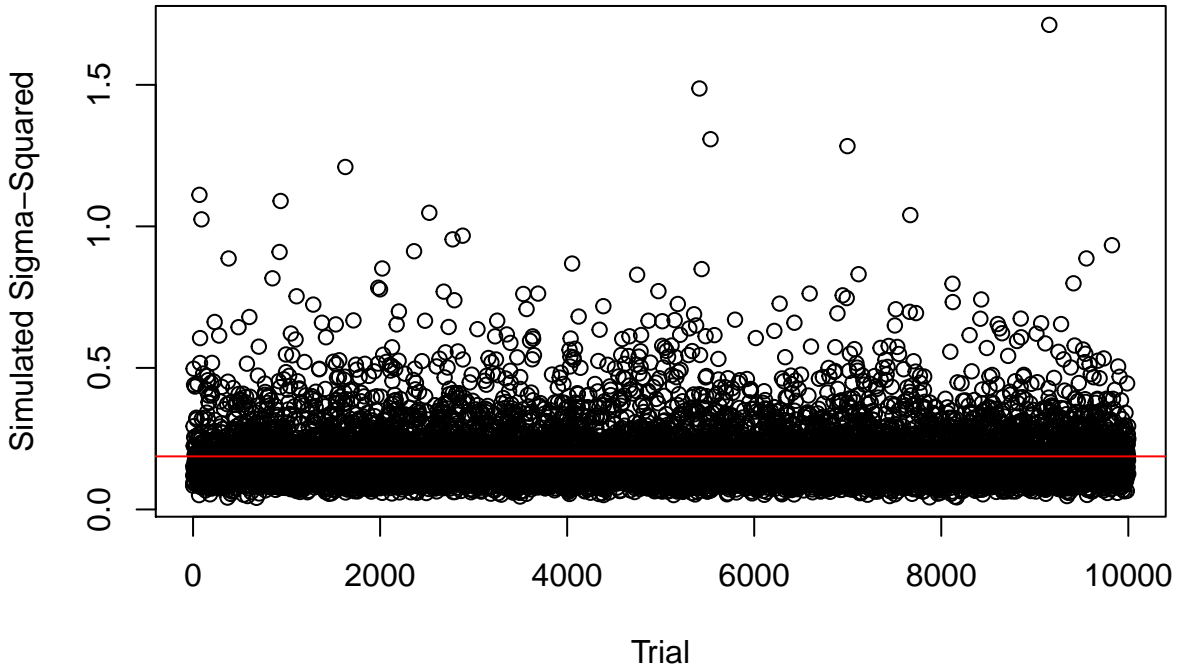
$$p(y|\mu, \sigma^2) = \frac{1}{y \cdot \sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2\sigma^2}(\log y - \mu)^2\right]$$

for $y > 0, \mu > 0, \sigma^2 > 0$ The log-normal distribution is related to the normal distribution as follows: if $y \sim \log N(\mu, \sigma^2)$ then $\log y \sim N(\mu, \sigma^2)$. Let $y_1, \dots, y_n | \mu, \sigma^2 \sim \log N(\mu, \sigma^2)$, where $\mu = 3.7$ is assumed to be known but σ^2 is unknown with non-informative prior $p(\sigma^2) \propto 1/\sigma^2$. The posterior for σ^2 is the $Inv - \chi^2(n, \tau^2)$ distribution, where

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

- (a) Simulate 10; 000 draws from the posterior of σ^2 (assuming $\mu = 3.7$) and compare with the theoretical $Inv - \chi^2(n, \tau^2)$ posterior distribution.

Comparing Simulated Values and Real Value

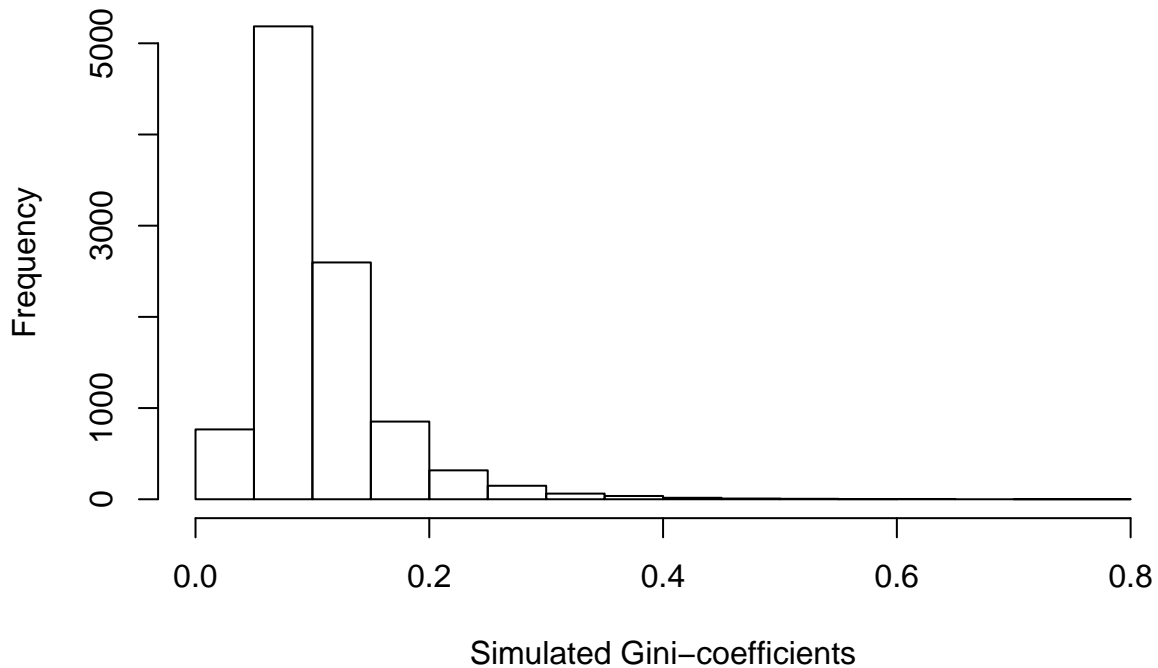


##	Simulation_Mean	Theoretical_Sigmasq	Difference_in_Percentage
## 1	0.1859347	0.1874264	0.795929

Difference between simulated mean and theoretical value in percentage is less than 1% ... very close values

- (b) 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. See Wikipedia for more information. It can be shown that $G = 2\Phi(\sigma/\sqrt{2}) - 1$ when incomes follow a $\log N(\mu, \sigma^2)$ distribution. $\Phi(z)$ is the cumulative distribution function (CDF) for the standard normal distribution with mean zero and unit variance. Use the posterior draws in a) to compute the posterior distribution of the Gini coefficient G for the current data set.

Histogram of Gini Coefficients from Simulation



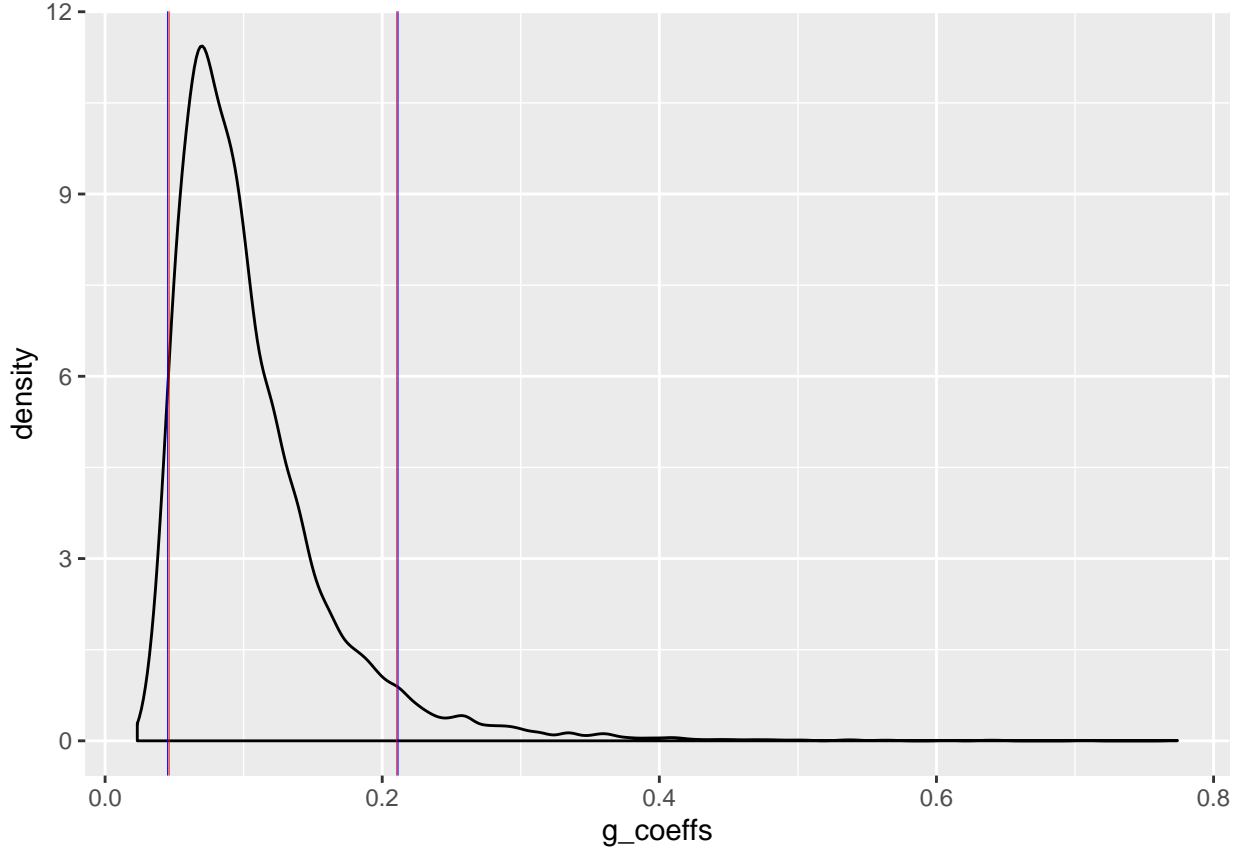
By generating 10000 draws and plotting the result, we can observe that the simulated Gini-coefficients are distributed around 0 to 0.4, which indicates this sample represents income inequality.

- (c) Use the posterior draws from b) to compute a 90% equal tail credible interval for G . A 90% equal tail interval (a ; b) cuts off 5% percent of the posterior probability mass to the left of a , and 5% to the right of b . Also, do a kernel density estimate of the posterior of G using the density function in R with default settings, and use that kernel density estimate to compute a 90% Highest Posterior Density interval for G . Compare the two intervals.

```
## 90% equal tail credible interval is:
```

```
## [1] 0.04618275 0.21053777
```

```
## 90% Highest Posterior Density interval is:
```



We can find that in this case, HPD has a wider interval comparing to ETI

3. Bayesian inference for the concentration parameter in the von Mises distribution. This exercise is concerned with directional data. The point is to show you that the posterior distribution for somewhat weird models can be obtained by plotting it over a grid of values. The data points are observed wind directions at a given location on ten different days. The data are recorded in degrees:

(40, 303, 326, 285, 296, 314, 20, 308, 299, 296)

where North is located at zero degrees (see Figure 1 on the next page, where the angles are measured clockwise). To fit with Wikipedias description of probability distributions for circular data we convert the data into radians $-\pi \leq y \leq \pi$ The 10 observations in radians are

(-2.44, 2.14, 2.54, 1.83, 2.02, 2.33, -2.79, 2.23, 2.07, 2.02)

Assume that these data points are independent observations following the von Mises distribution

$$p(y|\mu, k) = \frac{\exp[k \cdot \cos(y - \mu)]}{2\pi I_0(k)}, -\pi \leq y \leq \pi$$

where $I_0(k)$ is the modified Bessel function of the first kind of order zero [see ?besselI in R]. The parameter μ ($-\pi \leq y \leq \pi$) is the mean direction and $k > 0$ is called the concentration parameter. Large k gives a small variance around μ , and vice versa. Assume that μ is known to be 2:39. Let $k \sim \text{Exp}(\lambda = 1)$ a priori, where λ is the rate parameter of the exponential distribution (so that the mean is $1/\lambda$)

- (a) Plot the posterior distribution of k for the wind direction data over a fine grid of k values.

Likelihood for von Mises distribution

$$p(y|\mu, k) = \prod_{i=1}^n p(y_i|\mu, k)$$

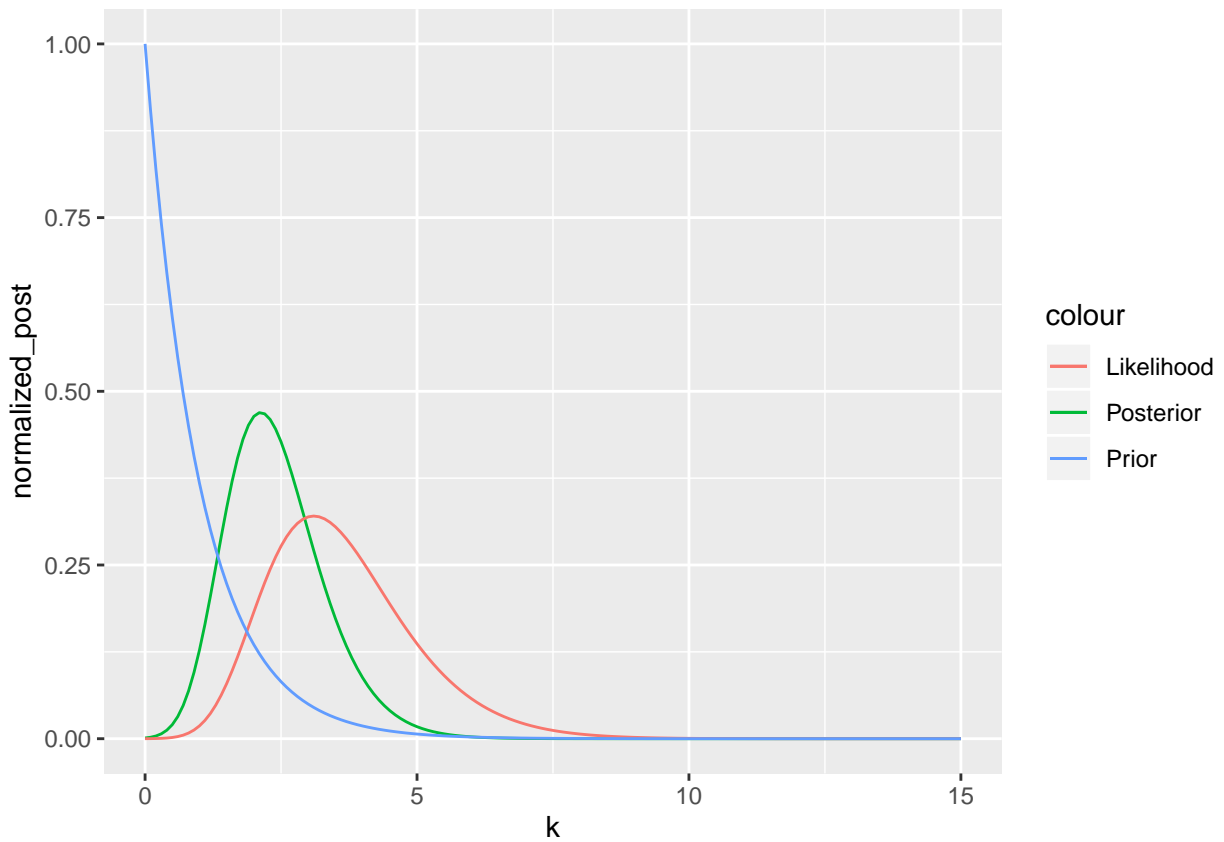
$$\begin{aligned}
&= \prod_{i=1}^n \frac{\exp[k \cdot \cos(y_i - \mu)]}{2\pi I_0(k)} \\
&= \left(\frac{1}{2\pi I_0(k)}\right)^n \cdot \exp\left\{\sum_{i=1}^n k \cdot \cos(y_i - \mu)\right\}
\end{aligned}$$

Prior

$$p(k) = \exp(\lambda = 1) = e^{-k}$$

Posterior = Likelihood x Prior

$$\begin{aligned}
p(k|y, u) &\propto p(y|\mu, k) \cdot p(k) \\
&\propto \frac{1}{I_0(k)} \cdot \exp\left(k \sum_{i=1}^n \cos(y_i - \mu)\right) \cdot \exp(-k) \\
&\propto \frac{1}{I_0(k)} \cdot \exp\left(k \sum_{i=1}^n \cos(y_i - \mu) - k\right)
\end{aligned}$$



After normalizing our prior, likelihood and posterior, we got the graph contains three of them with the area under the curve are 1.

(b) Find the (approximate) posterior mode of k from the information in a).

```
## [1] 3.1
```

Appendix A : Code Question 1

```
s=5
f=15
n=20
a0=2
beta0=2
true_a=a0+s
true_beta=beta0+f

means=c()
stds=c()

for(i in seq(10,10000,10)){
  pos=rbeta(i,true_a,true_beta)
  means=c(means,mean(pos))
  stds=c(stds,sd(pos))
}

true_mean=true_a/(true_a+true_beta)
true_var=(true_a*true_beta)/((true_a+true_beta+1)*(true_a+true_beta)^2)
cat("true mean=",true_mean)
cat("\ntrue standard deviation=",sqrt(true_var))
ggplot()+geom_point(aes(x=seq(10,10000,10),y=means))+geom_hline(yintercept=true_mean,col="red")+xlab("X")
ggplot()+geom_point(aes(x=seq(10,10000,10),y=stds))+geom_hline(yintercept = sqrt(true_var),col="red")+xlab("stds")
nDraws=10000
res=rbeta(nDraws,true_a,true_beta)
sample_prob=length(res[res>0.3])/nDraws
true_prob=1-pbeta(0.3,true_a,true_beta)
cat("sample prob=",sample_prob)
cat("\ntrue prob=",true_prob)
res=rbeta(nDraws,true_a,true_beta)
log_odds=log(pos/(1-pos))
hist(log_odds)
plot(density(log_odds))
```


Appendix B : Code Question 2

```
y = c(44,25,45,52, 30, 63, 19, 50, 34, 67)
log_y = log(y)
n = length(y)
mu = 3.7
tausq = sum((log_y - mu)^2)/n

sigmas = rinvchisq(n=10000, nu=n, tau = tausq)
plot(x=1:10000, y=sigmas, main="Comparing Simulated Values and Real Value", xlab="Trial", ylab="Simulated Values")
theoretical_val=n*tausq/(n-2)
abline(h=theoretical_val, col="red")
df=data.frame(Simulation_Mean=mean(sigmas), Theoretical_Sigmasq = theoretical_val, Difference_in_Percent=(Simulation_Mean-theoretical_val)/theoretical_val)
print(df)
g_coeffs = 2*pnorm(sigmas/sqrt(2)) -1
hist(g_coeffs, main="Histogram of Gini Coefficients from Simulation", xlab="Simulated Gini-coefficients")
g_coeffs_sorted = sort(g_coeffs)
a = g_coeffs_sorted[10000*0.05]
b = g_coeffs_sorted[10000*0.95]
cat("90% equal tail credible interval is:", "\n")
c(a,b)

# HPD
d=density(g_coeffs)
sim_dens=cumsum(d$y)/sum(d$y)
lower <- which(sim_dens>=0.05)[1]
upper <- which(sim_dens>=0.95)[1]
cat("90% Highest Posterior Density interval is:", "\n")
HPD=d$x[c(lower,upper)]

ggplot()+geom_density(aes(x=g_coeffs))+geom_vline(xintercept =a,col="red",size=0.01)+geom_vline(xintercept=b,col="red",size=0.01)
```

Appendix C : Code Question 3

```
mu=2.39
data=c(-2.44,2.14,2.54,1.83,2.02,2.33,-2.79,2.23,2.07,2.02)

k=seq(0,15,0.1)

posterior=function(k){
  return(exp(k*sum(cos(data-mu))-k)*besselI(k,0)^(-length(data)))
}
llik_fn=function(k){
  return(exp(k*sum(cos(data-mu)))*besselI(k,0)^(-length(data)))
}
prior_fn=function(k){
  return(exp(-k))
}
inte=integrate(posterior,lower = 0,upper = 15)
prior=prior_fn(k)
llik=exp(k*sum(cos(data-mu)))*(besselI(k,0))^(length(data))
normalized_llik=llik/integrate(llik_fn,0,15)$value
df=data.frame(k=k,normalized_post=(posterior(k)/inte$value),llik=normalized_llik,prior=prior)
ggplot(data=df)+geom_line(aes(x=k,y=normalized_post,colour="Posterior"))+geom_line(aes(x=k,y=llik,colour="Likelihood"))
#ggplot()+geom_line(aes(x=k,y=llik))
# let's choose the interval as 2 to 5
best_k_interval=abs(df$llik-df$normalized_post)[20:50]
best_k=which.min(best_k_interval)*0.1+2
best_k
```