

# Bayesian Learning - Lab 01

Lakshidaa Saigiridharan (laksa656) and Maximilian Pfundstein (marpf364)

2019-03-31

## Contents

<b>1</b>	<b>Bernoulli</b>	<b>1</b>
1.1	Drawing from the Posterior . . . . .	1
1.2	$\Pr(\theta < 0.4 y)$ . . . . .	4
1.3	Log-Odds . . . . .	4
<b>2</b>	<b>Log-Normal Distribution and the Gini Coefficient</b>	<b>5</b>
2.1	Simulate Draws from the Posterior Distribution . . . . .	6
2.2	Gini-Index . . . . .	7
2.3	Credible and Density Interval . . . . .	7
<b>3</b>	<b>Bayesian Inference</b>	<b>9</b>
3.1	Posterior Plot of $\kappa$ for Wind Direction Data Over a Fine Grid . . . . .	10
3.2	Posterior Mode of $\kappa$ . . . . .	11
<b>4</b>	<b>Source Code</b>	<b>12</b>

## 1 Bernoulli

**Exercise:** Let  $y_1, \dots, y_n | \theta \sim \text{Bern}(\theta)$ , 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 the posterior distribution  $\theta | y \sim \text{Beta}(\alpha_0 + s, \beta_0 + f)$ ,  $y = (y_1, \dots, y_n)$ , and verify that the posterior mean and standard deviation converges to the true values as the number of random draws grows large.
- Use simulation (`nDraws = 10000`) to compute the posterior propability  $\Pr(\theta < 0.4|y)$  and compare with the exact value [Hint: `pbeta()`].
- Compute the posterior distribution of the log-odds  $\phi = \log\left(\frac{\theta}{1-\theta}\right)$  by simulation (`nDraws = 10000`). [Hint: `hist()` and `density` might come in handy]

### 1.1 Drawing from the Posterior

First we define the parameters as we need them later.

```
#####  
# Exercise 1.a)  
#####  
  
# Parameters  
n = 20  
s = 14  
f = n - s
```

```

# Prior
alpha_z = 2
beta_z = 2

# Posterior
alpha_post = alpha_z + s
beta_post = beta_z + f

```

The posterior is given as  $\text{Beta}(\alpha_n, \beta_n)$  where  $\alpha_n = \alpha_0 + s$  and  $\beta_n = \beta_0 + f$ . Therefore the theoretical mean is given by:

$$E[X] = \frac{\alpha_n}{\alpha_n + \beta_n}$$

And the standard deviation by:

$$\text{sd}[X] = \sqrt{\frac{\alpha_n \beta_n}{(\alpha_n + \beta_n)^2 (\alpha_n + \beta_n + 1)}}$$

So let's calculate this.

```

mean_posterior = alpha_post / (alpha_post + beta_post)
sd_posterior = sqrt((alpha_post * beta_post) /
                    ((alpha_post + beta_post)^2 * (alpha_post + beta_post + 1)))

```

Therefore the mean of the prior is given by 0.6666667 and the standard deviation by 0.0942809.

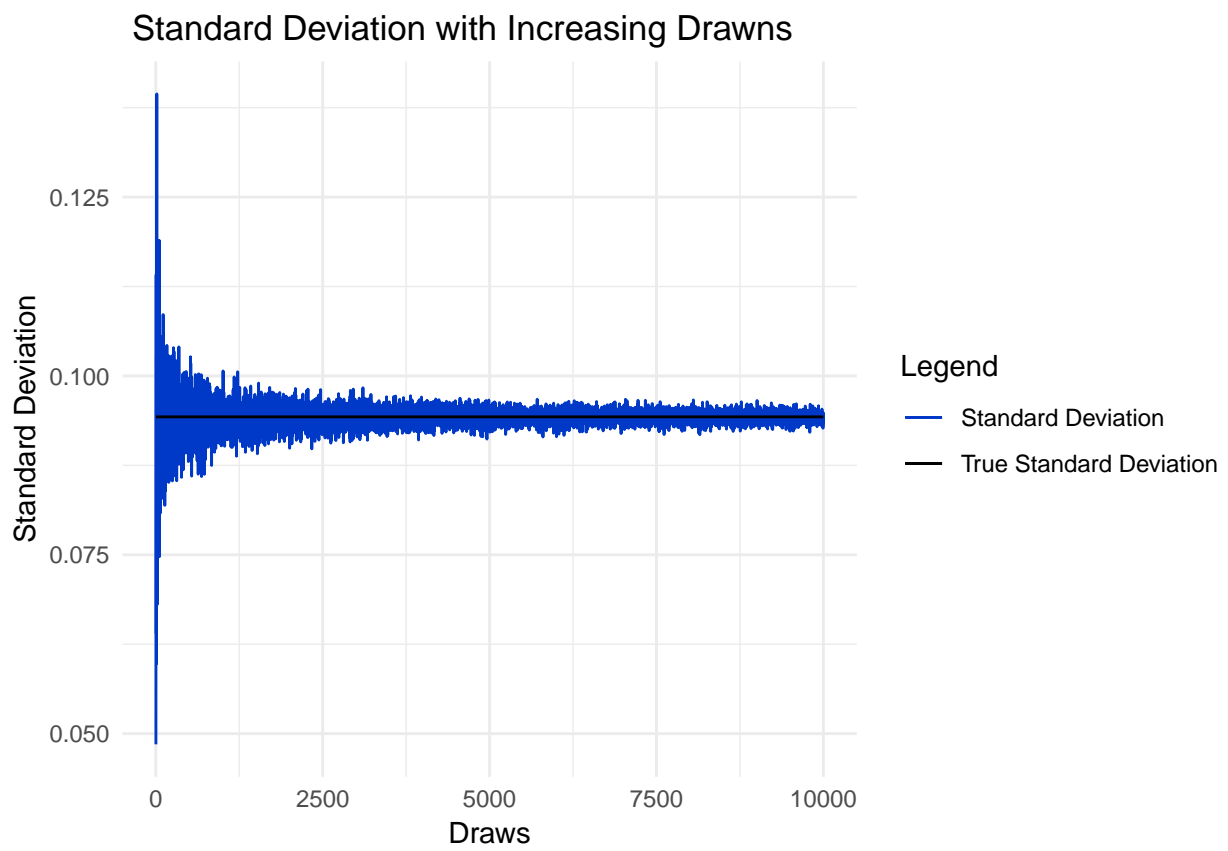
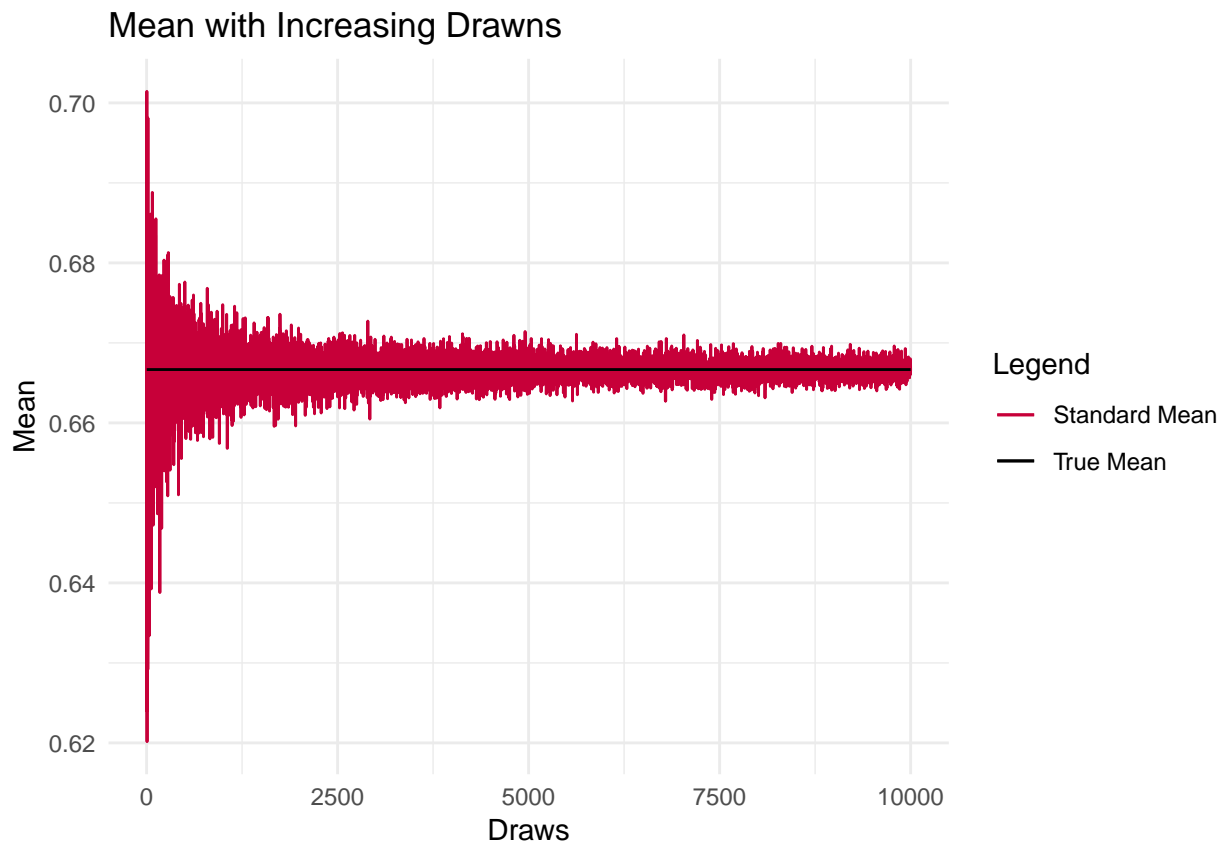
Now we will create a function that calculates the mean and standard deviation for a given number of trials to plot it later on.

```

get_stats = function(n, alpha, beta) {
  samples = rbeta(n, alpha, beta)
  return(c(count = n, sample_mean = mean(samples), sample_sd = sd(samples)))
}

df = data.frame(t(sapply(2:10000, get_stats, alpha_post, beta_post)))

```



We can see that the posterior mean and standard deviation converge as the number of random draws grows.

## 1.2 $Pr(\theta < 0.4|y)$

The true probability is calculated by calling `pbeta(0.4, alpha_post, beta_post)` which is 0.0039727.

We will simulate by taking samples and counting how many of them are  $< 0.4$ .

```
#####  
# Exercise 1.b)  
#####  
  
mean(rbeta(100000, alpha_post, beta_post) < 0.4)  
  
## [1] 0.00407
```

As we can see both values are quite close to each other.

## 1.3 Log-Odds

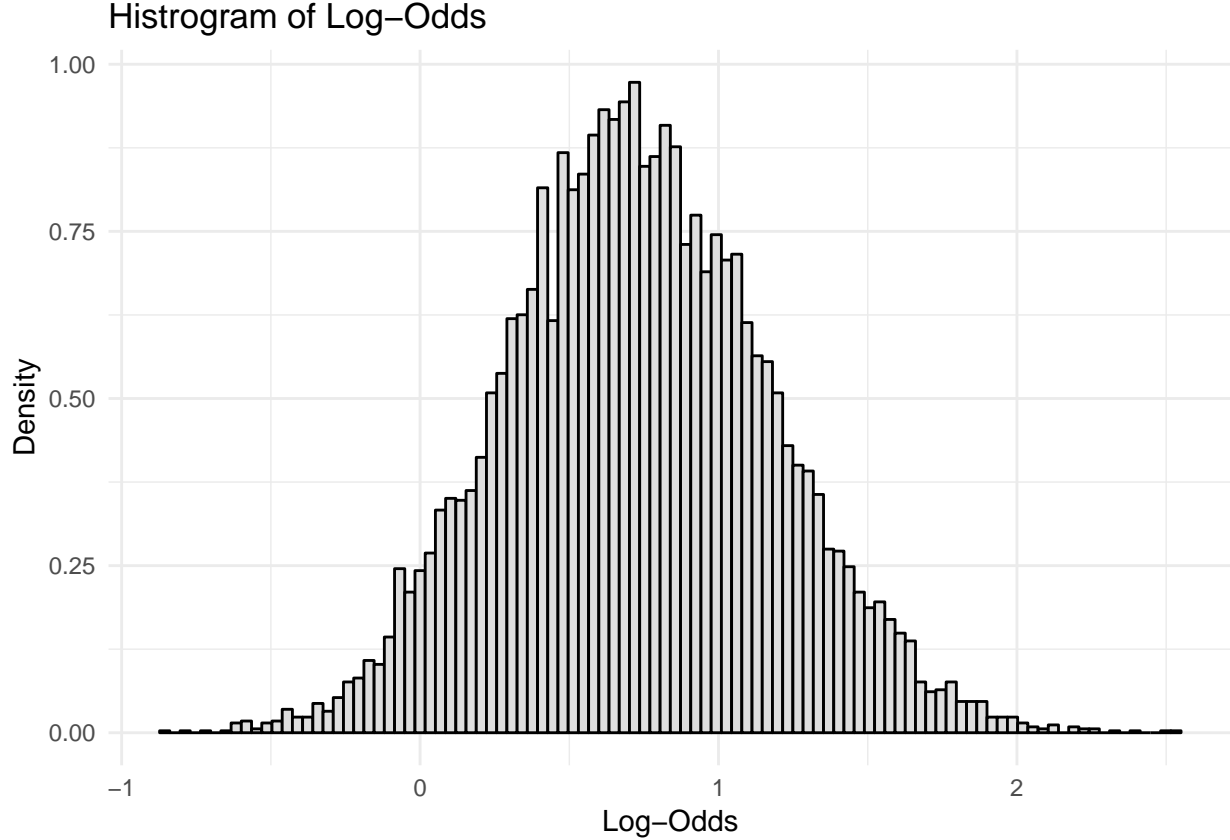
The log-odds are given by

$$\Phi = \log\left(\frac{\theta}{1-\theta}\right)$$

where  $\theta$  are samples drawn from the posterior. We can therefore easily calculate the value by:

```
#####  
# Exercise 1.c)  
#####  
  
draws = 10000  
  
samples = rbeta(draws, alpha_post, beta_post)  
phi = log(samples / (1 - samples))
```

The distribution looks as follows:



## 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: 14, 25, 45, 25, 30, 33, 19, 50, 34 and 67. A common model for non-negative continuous variables is the log-normal distribution. The log-normal distribution  $\mathcal{N}(\mu, \sigma^2)$  has density function

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

for  $y > 0$ ,  $\mu > 0$ , and  $\sigma^2 > 0$ . The log-normal distribution is related to the normal distribution as follows: if  $y \sim \mathcal{N}(\mu, \sigma^2)$  then  $\log y \sim \mathcal{N}(\mu, \sigma^2)$ . Let  $y_1, \dots, y_n | \mu, \sigma^2$ , where  $\mu = 3.5$  is assumed to be known but  $\sigma^2$  unknown with non-informative prior  $p(\sigma^2) \propto 1/\sigma^2$ . The posterior for  $\sigma^2$  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  $\sigma^2$  (assuming  $\mu = 3.5$ ) and compare with the theoretical  $Inv - \chi^2(n, \tau^2)$  posterior distribution.

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  $G = 2\Phi(\sigma/\sqrt{n}) - 1$  when incomes follow a  $\log \mathcal{N}(\mu, \sigma^2)$  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.

- c) Use the posterior draws from b) to compute a 95% equal tail credible interval for  $G$ . An 95% equal tail interval  $(a, b)$  cuts off 2.5% percent of the posterior probability mass to the left of  $a$ , and 97.5% 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 95% Highest Posterior Density interval for  $G$ . Compare the two intervals.

## 2.1 Simulate Draws from the Posterior Distribution

We calculate  $\tau^2$ , then we sample from `rinvchisq()` and compare the drawn samples to the density function given by `dinvchisq()`.

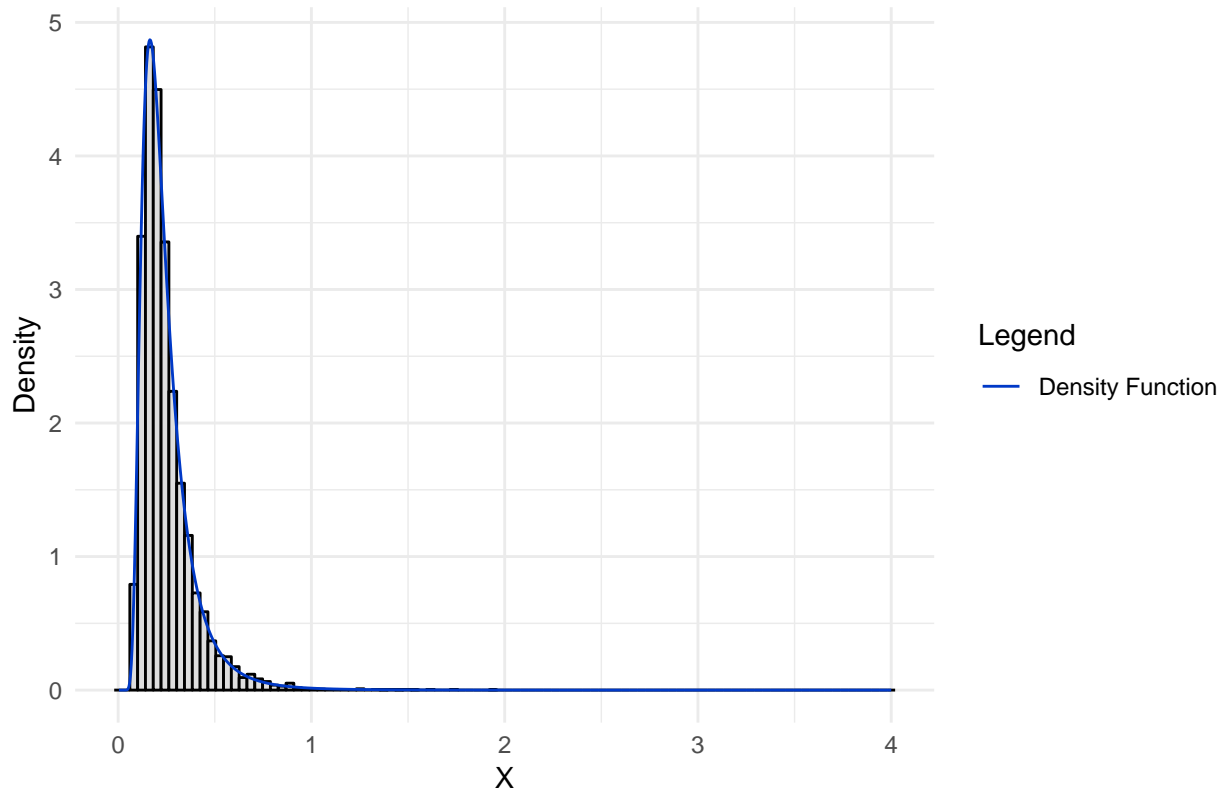
```
#####
# Exercise 2.a)
#####

obs = c(14, 25, 45, 25, 30, 33, 19, 50, 34, 67)
n = length(obs)
mu = 3.5
tau_sq = (sum((log(obs) - mu)^2)) / (n)

samples = rinvchisq(10000, n, tau_sq)
X = seq(from = 0, to = 4, length.out = 1000)
Y = dinvchisq(X, n, tau_sq)
```

We can see that the drawn samples for  $\sigma^2$  fit the  $Inv - \chi^2(n, \tau^2)$ .

### Drawn Samples and Inverse-Chi-Squared Density Function

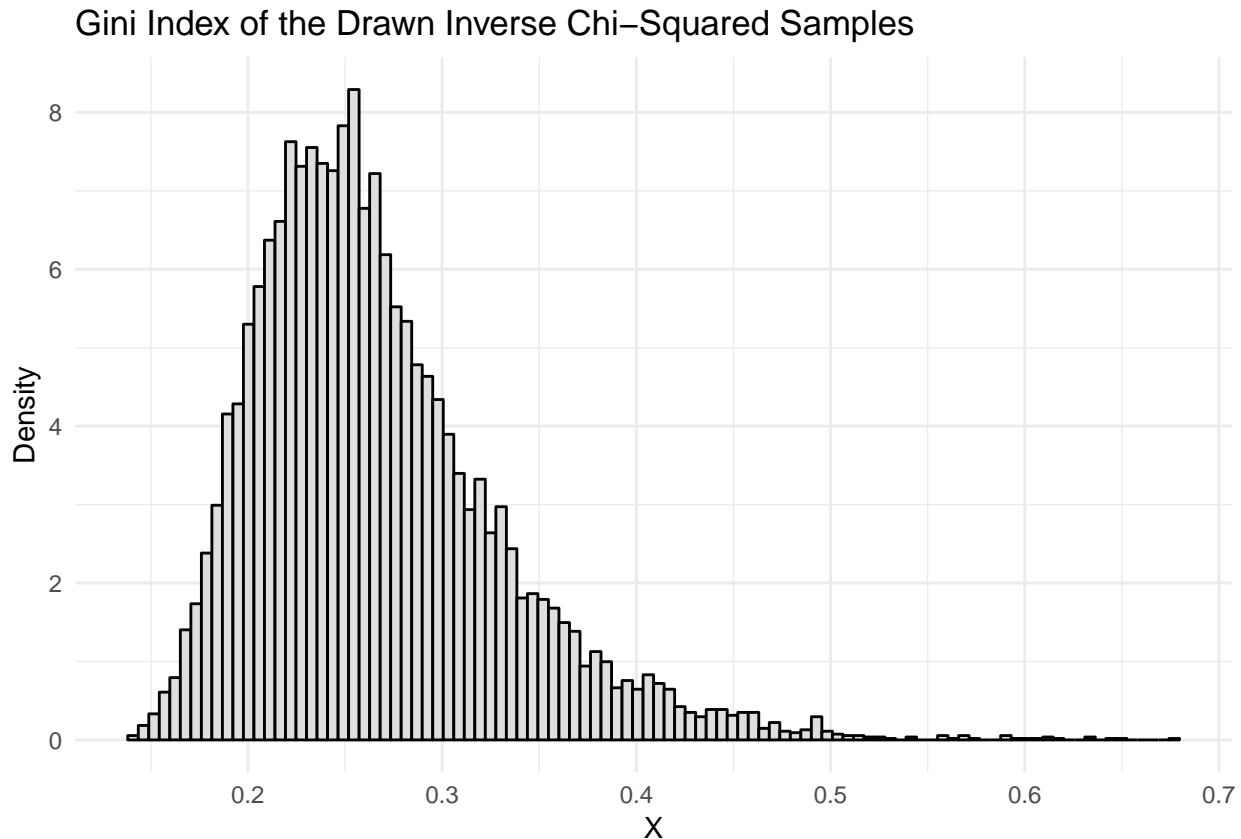


## 2.2 Gini-Index

The gini index is calculated by the following R code. As our samples represent  $\sigma^2$ , we have to include it in the square root.

```
#####  
# Exercise 2.b)  
#####  
  
G = 2 * pnorm(sqrt(samples/2)) - 1
```

The histogram for the gini indices looks like follows:



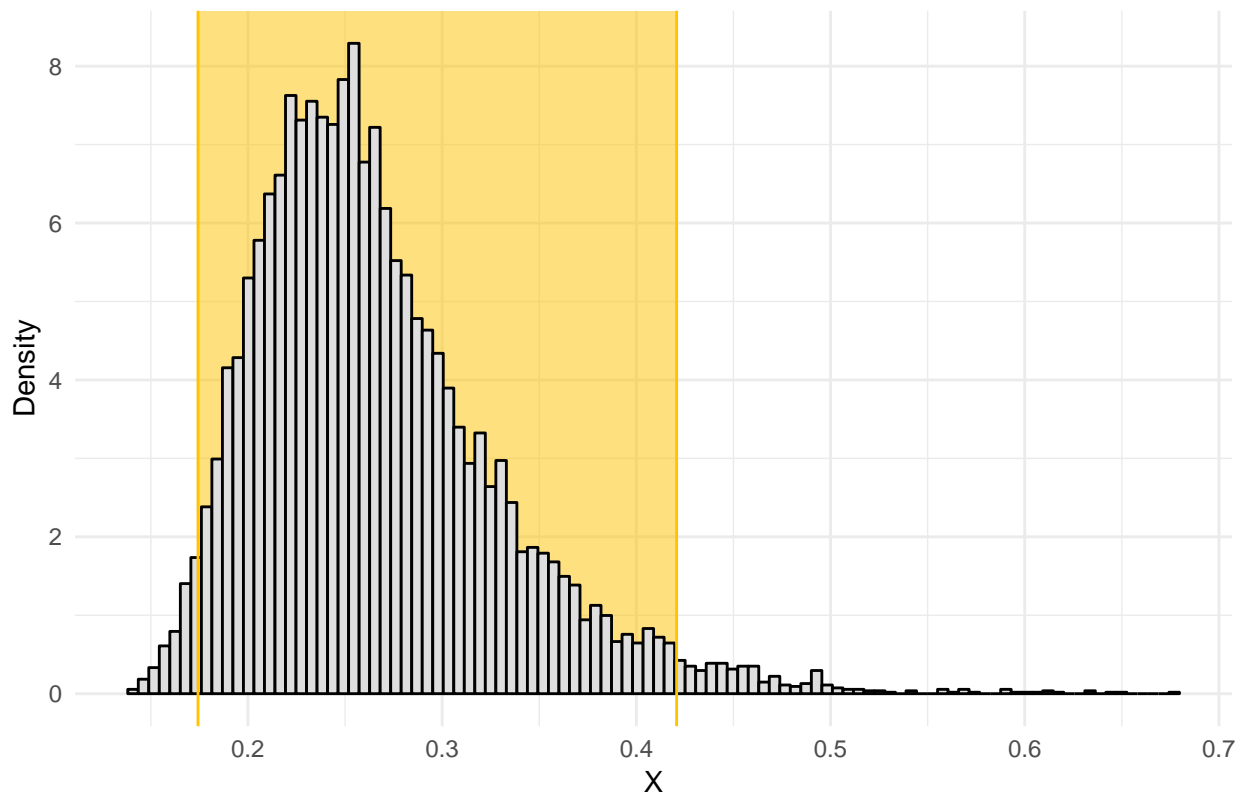
## 2.3 Credible and Density Interval

We use the `quantile()` function to receive the quantile values for the credible interval.

```
#####  
# Exercise 2.c)  
#####  
  
quantiles = quantile(G, c(0.025, 0.975))
```

The following plot visualizes the credible interval:

## Two Tailed 95% Credible Interval



For the density interval we first order our x and y values and then keep on adding more until we have reached more than 5 percent of the area (is faster then looking for the 95 percentile). We then save the index of that data point. The y value of this data point defines the threshold. We then filter for the values above the threshold, the left values all reside inside of the density interval. As we have just one interval in this case, it's enough to simply take the first and last value (of the unordered set after filtering) to visualize the density interval.

```
dg = density(G)

dg_unordered = data.frame(dg$x, dg$y)
dg_ordered = data.frame(dg$x[order(dg$y)], dg$y[order(dg$y)])
colnames(dg_ordered) = c("x_ord", "y_ord")

index = NaN

for (i in 1:nrow(dg_ordered)) {
  if (sum(dg_ordered$y_ord[1:i])/sum(dg_ordered$y_ord) > 0.05) {
    index = i
    break
  }
}

theta = dg_ordered[index,2]

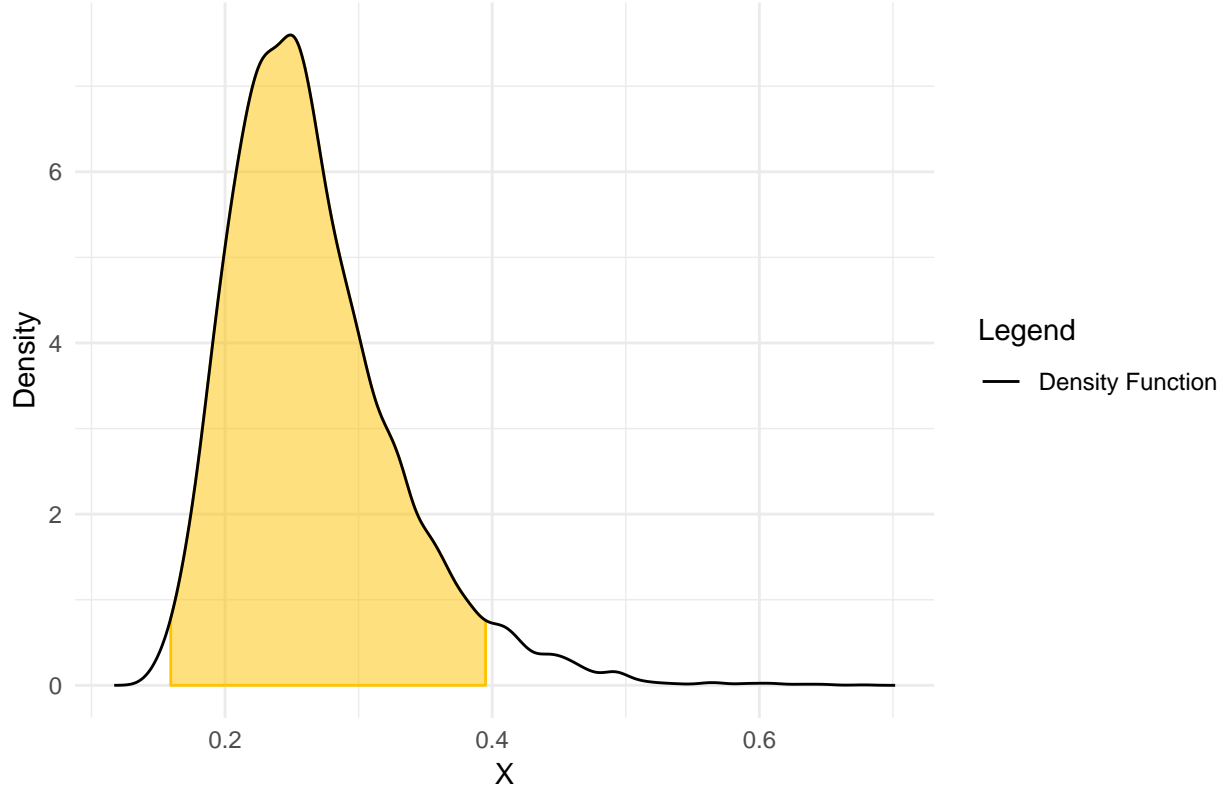
selected = dg_unordered[dg_unordered$dg.y > theta,]

interval_a = selected[1,]
interval_b = selected[nrow(selected),]
```



The value of  $\theta$  is 0.7493163. The interval is given by [0.1593518, 0.3950361].

### Gini Index of the Drawn Inverse Chi-Squared Samples



## 3 Bayesian Inference

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 Wikipedia's 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

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

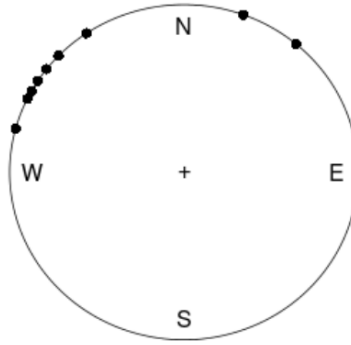


Figure 1: The wind direction data. Angles are measured clock-wise starting from North.

where  $I_0(\kappa)$  is the modified Bessel function of the first kind of order zero [see `?besselI` in R]. The parameter  $\mu(-\pi \leq y \leq \pi)$  is the mean direction and  $\kappa > 0$  is called the concentration parameter. Large  $\kappa$  gives a small variance around  $\mu$ , and vice versa. Assume that  $\mu$  is known to be 2.39. Let  $\kappa \sim \text{Exponential}(\lambda = 1)$  a priori, where  $\lambda$  is the rate parameter of the exponential distribution (so that the mean is  $1/\lambda$ ).

- Plot the posterior distribution of  $\kappa$  for the wind direction data over a fine grid of  $\kappa$  values.
- Find the (approximate) posterior mode of  $\kappa$  from the information in a).

### 3.1 Posterior Plot of $\kappa$ for Wind Direction Data Over a Fine Grid

```
#####
# Exercise 3.a)
#####

mu = 2.39
lambda = 1
kappa = seq(from = 0.01, to = 20.0, by = 0.01)
y = c(-2.44, 2.14, 2.54, 1.83, 2.02, 2.33, -2.79, 2.23, 2.07, 2.02)

p = function(y, kappa, mu = 2.39) {
  return((exp(kappa * cos(y - mu))) / (2 * pi * besselI(kappa, 0)))
}

likelihood = function(y, kappa, mu = 2.39) {
  m_probs = sapply(kappa, FUN = function(k, y, mu) {
    return(p(y, k, mu))
  }, y, mu)
  return(apply(m_probs, 2, prod))
}

prior = function(kappa, lambda) {
  return(dexp(kappa, rate = lambda))
}
```

```

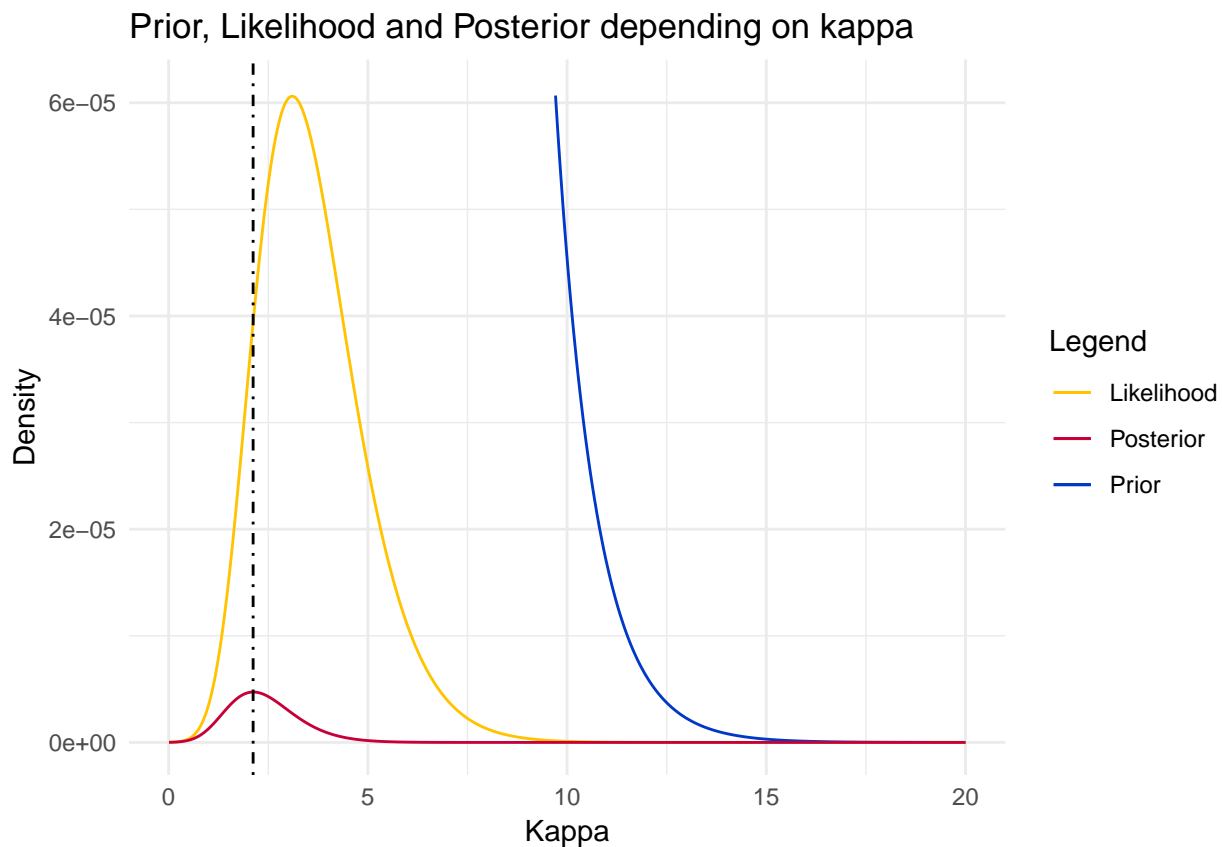
posterior = function(y, kappa, mu = 2.39, lambda = 1) {
  return(prior(kappa, lambda) * likelihood(y, kappa, mu))
}

prior_values = prior(kappa, lambda)
likelihood_values = likelihood(y, kappa, mu)
# posterior(y, kappa, mu) would work as well but as we already have the values
# we can save some computation time
posterior_values = prior_values * likelihood_values

#####
# Exercise 3.b)
#####

index = which.max(posterior_values)
kappa_max = kappa[index]

```



### 3.2 Posterior Mode of $\kappa$

For approximating the posterior mode we have to find the maximum value of the posterior and find the appropriate  $\kappa$ . The code can be found in the code snippet above. The value is marked by the dashed line in the plot.

The posterior mode is therefore given by  $\kappa = 2.12$ .

## 4 Source Code

```
knitr::opts_chunk$set(echo = TRUE)
library(ggplot2)
library(geoR)

#####
# Exercise 1.a)
#####

# Parameters
n = 20
s = 14
f = n - s

# Prior
alpha_z = 2
beta_z = 2

# Posterior
alpha_post = alpha_z + s
beta_post = beta_z + f

mean_posterior = alpha_post / (alpha_post + beta_post)
sd_posterior = sqrt((alpha_post * beta_post) /
                    ((alpha_post + beta_post)^2 * (alpha_post + beta_post + 1)))

get_stats = function(n, alpha, beta) {
  samples = rbeta(n, alpha, beta)
  return(c(count = n, sample_mean = mean(samples), sample_sd = sd(samples)))
}

df = data.frame(t(apply(2:10000, get_stats, alpha_post, beta_post)))

ggplot(df) +
  geom_line(aes(x = count, y = sample_mean, color = "Standard Mean")) +
  geom_line(aes(x = count, y = mean_posterior, color = "True Mean")) +
  labs(title = "Mean with Increasing Draws", y = "Mean", x = "Draws") +
  scale_color_manual("Legend", values = c("#C70039", "#000000")) +
  theme_minimal()

ggplot(df) +
  geom_line(aes(x = count, y = sample_sd, colour = "Standard Deviation")) +
  geom_line(aes(x = count, y = sd_posterior, colour = "True Standard Deviation")) +
  labs(title = "Standard Deviation with Increasing Draws",
       y = "Standard Deviation", x = "Draws") +
  scale_color_manual("Legend", values = c("#0039C7", "#000000")) +
  theme_minimal()
```

```
#####
# Exercise 1.b)
#####

mean(rbeta(100000, alpha_post, beta_post) < 0.4)

#####
# Exercise 1.c)
#####

draws = 10000

samples = rbeta(draws, alpha_post, beta_post)
phi = log(samples / (1 - samples))

phi = data.frame(phi)
colnames(phi) = "phi"

ggplot(phi) +
  geom_histogram(aes(x = phi, y = ..density..), color = "black",
                 fill = "#dedede", bins = sqrt(draws)) +
  labs(title = "Histogram of Log-Odds",
       y = "Density",
       x = "Log-Odds", color = "Legend") +
  theme_minimal()

#####
# Exercise 2.a)
#####

obs = c(14, 25, 45, 25, 30, 33, 19, 50, 34, 67)
n = length(obs)
mu = 3.5
tau_sq = (sum((log(obs) - mu)^2)) / (n)

samples = rinvchisq(10000, n, tau_sq)
X = seq(from = 0, to = 4, length.out = 1000)
Y = dinvchisq(X, n, tau_sq)

df = data.frame(X, Y)
samples_df = data.frame(samples)

ggplot(df) +
  geom_histogram(data = samples_df, aes(x = samples, y = ..density..),
                 bins = sqrt(nrow(samples_df)), color = "black",
                 fill = "#DEDEDE") +
  geom_line(aes(x = X, y = Y, color = "Density Function")) +
  labs(title = "Drawn Samples and Inverse-Chi-Squared Density Function",
       y = "Density", x = "X") +
```

```

scale_color_manual("Legend", values = c("#0039C7", "#000000")) +
theme_minimal()

#####
# Exercise 2.b)
#####

G = 2 * pnorm(sqrt(samples/2)) - 1

G_df = data.frame(G)

ggplot(G_df) +
  geom_histogram(aes(x = G, y=..density..),
                 bins = sqrt(nrow(G_df)), color = "black", fill = "#DEDEDE") +
  labs(title = "Gini Index of the Drawn Inverse Chi-Squared Samples",
       y = "Density", x = "X") +
  scale_color_manual("Legend", values = c("#0039C7", "#000000")) +
  theme_minimal()

#####
# Exercise 2.c)
#####

quantiles = quantile(G, c(0.025, 0.975))

ggplot(G_df) +
  annotate("rect", xmin=quantiles[1], xmax=quantiles[2], ymin=0, ymax=Inf,
         alpha=0.5, fill="#FFC300") +
  geom_histogram(aes(x = G, y=..density..),
                 bins = sqrt(nrow(G_df)), color = "black", fill = "#DEDEDE") +
  geom_vline(xintercept = quantiles, colour = "#FFC300") +
  labs(title = "Two Tailed 95% Credible Interval",
       y = "Density", x = "X") +
  scale_color_manual("Legend", values = c("#0039C7", "#000000")) +
  theme_minimal()

dg = density(G)

dg_unordered = data.frame(dg$x, dg$y)
dg_ordered = data.frame(dg$x[order(dg$y)], dg$y[order(dg$y)])
colnames(dg_ordered) = c("x_ord", "y_ord")

index = NaN

for (i in 1:nrow(dg_ordered)) {
  if (sum(dg_ordered$y_ord[1:i])/sum(dg_ordered$y_ord) > 0.05) {
    index = i
    break
  }
}

```

```

}

theta = dg_ordered[index,2]

selected = dg_unordered[dg_unordered$dg.y > theta,]

interval_a = selected[1,]
interval_b = selected[nrow(selected),]

dg_df = data.frame(dg$x, dg$y)

ggplot(dg_df) +
  geom_ribbon(data = selected, aes(x = dg.x, ymin = 0, ymax = dg.y),
            alpha = 0.5, fill = "#FFC300", color = "#FFC300") +
  geom_line(aes(x = dg.x, y = dg.y, color = "Density Function")) +
  labs(title = "Gini Index of the Drawn Inverse Chi-Squared Samples",
       y = "Density", x = "X") +
  scale_color_manual("Legend", values = c("#000000", "#000000")) +
  theme_minimal()

#####
# Exercise 3.a)
#####

mu = 2.39
lambda = 1
kappa = seq(from = 0.01, to = 20.0, by = 0.01)
y = c(-2.44, 2.14, 2.54, 1.83, 2.02, 2.33, -2.79, 2.23, 2.07, 2.02)

p = function(y, kappa, mu = 2.39) {
  return((exp(kappa * cos(y - mu))) / (2 * pi * besseli(kappa, 0)))
}

likelihood = function(y, kappa, mu = 2.39) {
  m_probs = sapply(kappa, FUN = function(k, y, mu) {
    return(p(y, k, mu))
  }, y, mu)
  return(apply(m_probs, 2, prod))
}

prior = function(kappa, lambda) {
  return(dexp(kappa, rate = lambda))
}

posterior = function(y, kappa, mu = 2.39, lambda = 1) {
  return(prior(kappa, lambda) * likelihood(y, kappa, mu))
}

prior_values = prior(kappa, lambda)
likelihood_values = likelihood(y, kappa, mu)
# posterior(y, kappa, mu) would work as well but as we already have the values

```

```

# we can save some computation time
posterior_values = prior_values * likelihood_values

#####
# Exercise 3.b)
#####

index = which.max(posterior_values)
kappa_max = kappa[index]

df = data.frame(kappa, prior_values, likelihood_values, posterior_values)
colnames(df) = c("kappa", "prior", "likelihood", "posterior")

ggplot(df) +
  geom_line(aes(x = kappa, y = prior, color = "Prior")) +
  geom_line(aes(x = kappa, y = likelihood, color = "Likelihood")) +
  geom_line(aes(x = kappa, y = posterior, color = "Posterior")) +
  geom_vline(aes(xintercept = kappa_max), linetype = "dotted") +
  labs(title = "Prior, Likelihood and Posterior depending on kappa",
       y = "Density", x = "Kappa") +
  scale_color_manual("Legend", values = c("#FFC300", "#C70039", "#0039C7")) +
  ylim(0, 6.1*10-5) +
  theme_minimal()

```