

# STAT 656: Bayesian Data Analysis

## Fall 2024

### Homework 1

Juanwu Lu\*

## Synthetic Data

The *autoregressive model* is frequently used to analyze time series data. The simplest autoregressive model has order 1, and is abbreviated as AR(1). This model assumes that an observation  $y_i$  at time point  $i$  ( $i = 1, \dots, n$ ) is generated according to

$$y_i = \rho y_{i-1} + \epsilon_i,$$

where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$  independently, and  $\rho$  and  $\sigma$  are unknown parameters. For simplicity, we shall assume that  $y_0$  is a fixed constant. We will also assume  $|\rho| < 1$ .

1. (5 points) **Solution:**

Given the formulation above, the log-likelihood function is calculated as follows:

$$\begin{aligned} \log L(\rho, \sigma^2 | y_0, y_1, \dots, y_n) &= \log \prod_{i=1}^n (2\pi\sigma^2)^{-\frac{1}{2}} \cdot \exp \left\{ -\frac{(y_i - \rho y_{i-1})^2}{2\sigma^2} \right\} \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \rho y_{i-1})^2 \\ &= -\frac{n}{2} \log(2\pi) - n \log(\sigma) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \rho y_{i-1})^2. \end{aligned}$$

2. (10 points) **Solution:**

The code implementation is as follows:

```
# Read data
if (file.exists("data/computation_data_hw_1.csv")) {
  data <- read.csv("data/computation_data_hw_1.csv")
  y <- data[["x"]]
} else {
  stop("Cannot find the file 'data/computation_data_hw_1.csv' at ", getwd())
}

# Define the log-likelihood function
ar_loglik <- function(rho, log_sig) {
  y <- .GlobalEnv$y
  n <- length(y)
  sig <- exp(log_sig)
```

---

\*College of Engineering, Purdue University, West Lafayette, IN, USA

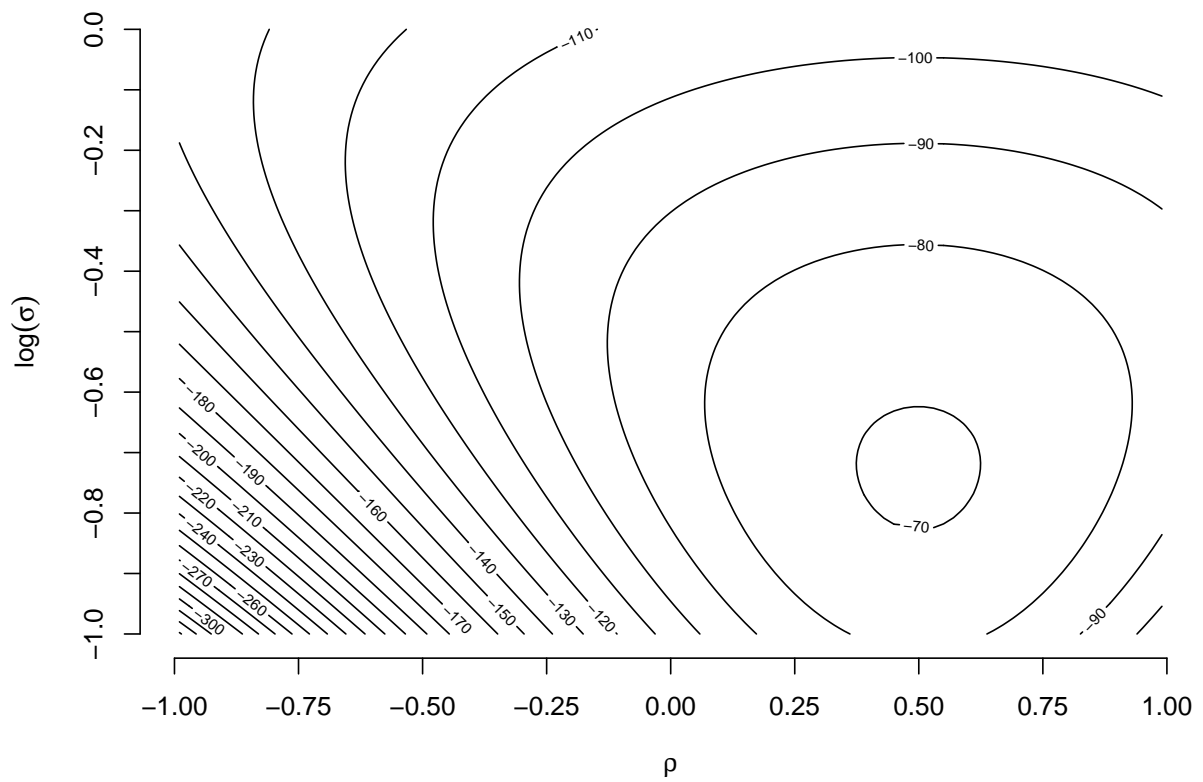
```

rho <- rep(as.numeric(rho), times = n - 1)
denom <- n * log(2 * pi) + 2 * n * log_sig
log_lik <- -0.5 * (
  denom + (y[1]^2 + sum((y[2:n] - rho * y[1:(n - 1)])^2)) / sig^2
)
return(log_lik)
}

# Visualization
rho <- seq(-0.99, 0.99, length = 100)
log_sig <- seq(-1.0, 0.0, length = 100)
loglik <- outer(rho, log_sig, Vectorize(ar_loglik))
contour(
  x = rho,
  y = log_sig,
  z = loglik,
  xlab = expression(rho),
  ylab = expression(log(sigma)),
  main = "Log-Likelihood Function",
  nlevels = 20,
  axes = FALSE,
)
axis(side = 1, at = seq(-1.0, 1.0, by = 0.25))
axis(side = 2, at = seq(-1.0, 0.0, by = 0.10))

```

**Log-Likelihood Function**



3. (10 points) **Solution:**

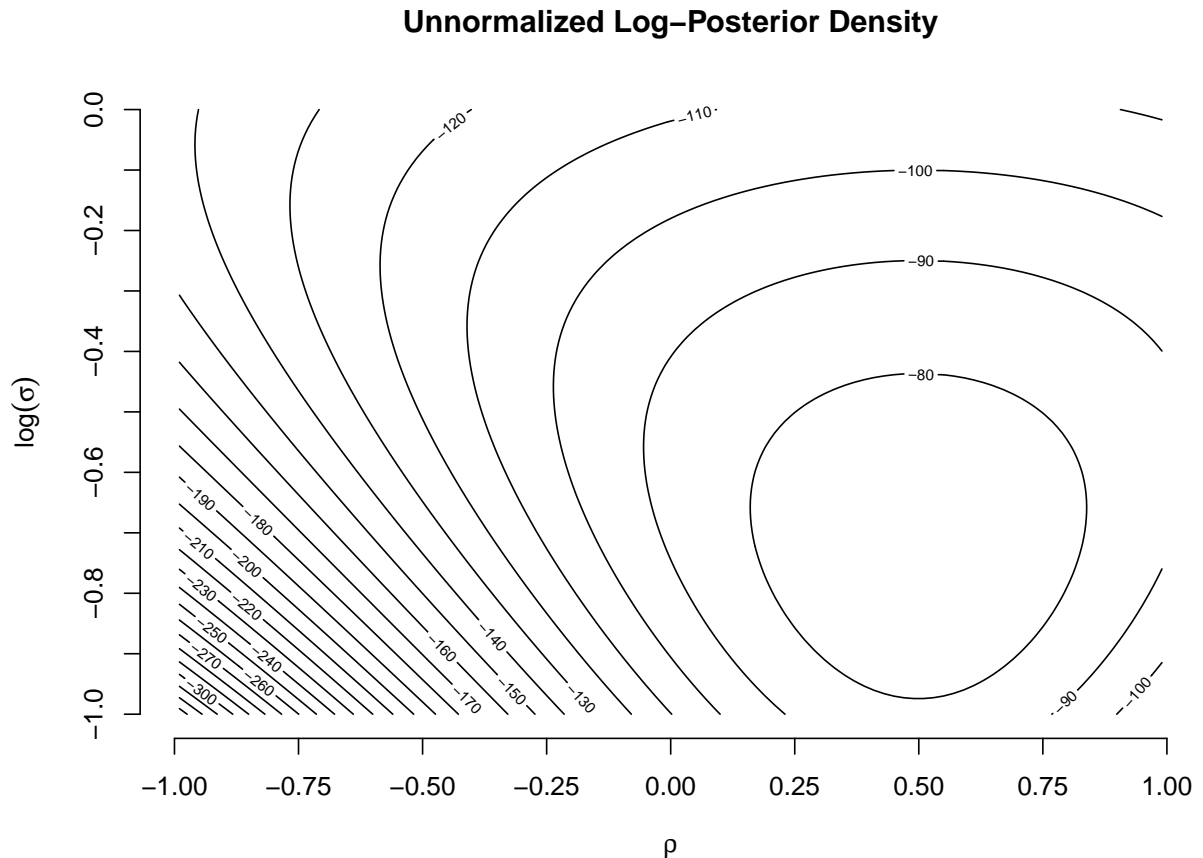
Since the prior is independent, we have  $p(\rho, \log(\sigma)) = p(\rho)p(\log(\sigma)) = \frac{1}{2} \cdot \frac{1}{10\sqrt{2\pi}} \exp\left\{-\frac{\log(\sigma)^2}{200}\right\}$ . Therefore, the posterior density is proportional to the product of the likelihood and the prior, and then the log of the posterior density is calculated up to a constant as follows:

$$\begin{aligned}\log p(\rho, \log(\sigma)|y_0, y_1, \dots, y_n) &= \log L(\rho, \log(\sigma)|y_0, y_1, \dots, y_n) + \log p(\rho, \log(\sigma)) + C \\ &= -n \log(\sigma) - \frac{1}{2 \exp(2 * \log(\sigma))} \sum_{i=1}^n (y_i - \rho y_{i-1})^2 - \frac{\log(\sigma)^2}{200} + C'\end{aligned}$$

where  $C$  is the constant log-normalizer and  $C'$  is a constant irrelevant to  $\rho$  and  $\log(\sigma)$ . The visualization of the log of the posterior density is as follows:

```
ar_logpost <- function(rho, log_sig) {
  log_prior <- -log(2) - 0.5 * (log(2 * pi) + log(100) + log_sig^2 / 100)
  log_post <- ar_loglik(rho, log_sig) + log_prior
  return(log_post)
}

logpost <- outer(rho, log_sig, Vectorize(ar_logpost))
contour(
  x = rho,
  y = log_sig,
  z = logpost,
  main = "Unnormalized Log-Posterior Density",
  xlab = expression(rho),
  ylab = expression(log(sigma)),
  nlevels = 20,
  axes = FALSE,
)
axis(side = 1, at = seq(-1.0, 1.0, by = 0.25))
axis(side = 2, at = seq(-1.0, 0.0, by = 0.10))
```



Compared to the log-likelihood function, the log posterior density is more concentrated around the maximum likelihood estimate of  $(\rho, \log(\sigma))^T$ . The prior is not overly informative since the shape of the posterior density is still similar to the likelihood function, indicating that the posterior is **mainly determined by the likelihood function**.

4. (10 points) **Solution:**

From the visualization in 3., we can see that the posterior density is concentrated around  $0.25 \leq \rho \leq 0.75$  and  $-0.9 \leq \log(\sigma) \leq -0.5$ . Therefore, we can choose a grid of  $\rho$  and  $\log(\sigma)$  as follows:

```
set.seed(42)
rho_grid <- seq(0.25, 0.75, length = 100)
log_sig_grid <- seq(-0.90, -0.50, length = 100)
grid_log_post <- outer(rho_grid, log_sig_grid, Vectorize(ar_logpost))
# Calculate the normalized probability density
probs <- exp(grid_log_post - max(grid_log_post))
probs <- probs / sum(probs)
# Randomly sample from the grid
indices <- sample(
  x = seq_len(length(as.vector(probs))),
  size = 1000,
  replace = TRUE,
  prob = probs
)
rho_sample <- rho_grid[((indices - 1) %% nrow(probs)) + 1]
log_sig_sample <- log_sig_grid[((indices - 1) %% nrow(probs)) + 1]
```

5. (5 points) **Solution:**

The code implementation is as follows:

```
library(moments)

# Calculate summaries for rho
print(quantile(rho_sample, probs = c(0.025, 0.25, 0.5, 0.75, 0.975)))

##      2.5%      25%      50%      75%      97.5%
## 0.3357323 0.4419192 0.4974747 0.5542929 0.6691919

sprintf("Mean: %.4f", mean(rho_sample))

## [1] "Mean: 0.5002"

sprintf("Standard Deviation: %.4f", sd(rho_sample))

## [1] "Standard Deviation: 0.0847"

sprintf("Skewnewss: %.4f", skewness(rho_sample))

## [1] "Skewnewss: 0.0490"

sprintf("Kurtosis: %.4f", kurtosis(rho_sample))

## [1] "Kurtosis: 3.0670"

# Calculate summaries for log(sigma)
print(quantile(log_sig_sample, probs = c(0.025, 0.25, 0.5, 0.75, 0.975)))

##      2.5%      25%      50%      75%      97.5%
## -0.8516162 -0.7707071 -0.7222222 -0.6696970 -0.5646465

sprintf("Mean: %.4f", mean(log_sig_sample))

## [1] "Mean: -0.7199"

sprintf("Standard Deviation: %.4f", sd(log_sig_sample))

## [1] "Standard Deviation: 0.0710"

sprintf("Skewness: %.4f", skewness(log_sig_sample))

## [1] "Skewness: 0.1683"

sprintf("Kurtosis: %.4f", kurtosis(log_sig_sample))

## [1] "Kurtosis: 2.7730"
```

6. (10 points) **Solution:**

The code implementation is as follows:

```
ar_post_predictive <- function(rho, log_sig) {
  # Sample from the posterior predictive distribution
  y <- .GlobalEnv$y
  new_y <- rep(0.0, times = length(y))
  sig <- exp(log_sig)
  for (i in 2:length(y)) {
    new_y[i] <- rnorm(n = 1, mean = rho * y[i - 1], sd = sig)
  }
  return(new_y)
}
```

```

params <- data.frame(rho = rho_sample, log_sig = log_sig_sample)
samples <- matrix(NA, nrow = nrow(params), ncol = length(y))
for (i in seq_len(nrow(params))) {
  samples[i, ] <- ar_post_predictive(params[i, "rho"], params[i, "log_sig"])
}

```

The above code makes use of the grid samples drawn in problem 4. For each pair of  $(\rho, \log(\sigma))$ , we use the AR(1) model to generate a new sample of sequences. The summary statistics of the posterior predictive distribution are as follows:

```
print("Sequence means:")
```

```
## [1] "Sequence means:"
```

```
print(colMeans(samples))
```

```

##      [1]  0.000000000 -0.004768199  0.380467600  0.046329176  0.041871692
##      [6]  0.230962145 -0.239262815  0.214538988  0.320248721  0.833946396
##     [11] -0.047636240  0.154295216  0.091970943  0.491005558  0.167442727
##     [16]  0.448882392  0.221520741  0.374755026 -0.241297744 -0.400832166
##     [21]  0.027256619 -0.401138975  0.116944379 -0.011085951  0.170830053
##     [26] -0.044556151 -0.312264947 -0.387679056 -0.729715619 -0.668030634
##     [31] -0.604599248 -0.828631029 -0.414124777 -0.153571316 -0.065036616
##     [36]  0.016159584  0.160684738  0.128898047 -0.165708086  0.295168094
##     [41]  0.400073542 -0.100819868  0.057687704 -0.194070728 -0.274464118
##     [46] -0.054924127 -0.393364772 -0.596120509 -0.254331298  0.058619835
##     [51]  0.161324834  0.058085082 -0.074207172  0.054033284  0.271008982
##     [56]  0.005264908 -0.169212597  0.181303334  0.197683160  0.174879318
##     [61]  0.409731341 -0.031243178 -0.011589809  0.019277832  0.164153831
##     [66] -0.070950225 -0.322370283 -0.174962173 -0.134343046 -0.213333649
##     [71] -0.360695027  0.253886889 -0.063299359 -0.218656995 -0.072649060
##     [76]  0.038109116  0.091598943  0.439842892  0.175992690 -0.205825087
##     [81]  0.249317439 -0.038973702  0.255908060 -0.029138104 -0.056703887
##     [86] -0.183220258 -0.158864768 -0.146437883 -0.186715701 -0.252234505
##     [91] -0.172413509  0.005136642 -0.111024282  0.241944648  0.249826061
##     [96]  0.067545792  0.219785394  0.366799730  0.354160286 -0.187007079

```

```
print("Sequence standard deviations:")
```

```
## [1] "Sequence standard deviations:"
```

```
print(apply(samples, 2, sd))
```

```

##      [1] 0.0000000 0.4875799 0.4952396 0.4853952 0.4985296 0.5092669 0.4943722
##      [8] 0.5042714 0.4834231 0.4975461 0.4991716 0.4769155 0.4994349 0.4932221
##     [15] 0.4738587 0.4824026 0.4765846 0.4818143 0.5206273 0.4930490 0.4913379
##     [22] 0.4697823 0.5076087 0.4998257 0.4994883 0.4776279 0.5204864 0.4897410
##     [29] 0.5075770 0.5083453 0.4970612 0.5102311 0.5127232 0.4965778 0.4776371
##     [36] 0.4897456 0.4745159 0.4765058 0.4915986 0.4936407 0.4896082 0.5106595
##     [43] 0.4947319 0.4906189 0.5111469 0.4905928 0.4780761 0.5285908 0.5027156
##     [50] 0.4952306 0.4885762 0.4887156 0.4954808 0.4826266 0.5065151 0.4840795
##     [57] 0.4783016 0.4926563 0.5074513 0.4827778 0.4915453 0.4948413 0.4831472
##     [64] 0.4882579 0.4893974 0.4978840 0.4923274 0.4964880 0.4812996 0.5041211
##     [71] 0.4991301 0.4958067 0.5055281 0.4947765 0.4903025 0.4971520 0.4988327
##     [78] 0.4844680 0.4888304 0.4860582 0.4848816 0.4869918 0.4963936 0.4814153
##     [85] 0.4818288 0.5007236 0.4806716 0.5092220 0.4856871 0.4953563 0.4998521
##     [92] 0.4754400 0.4945970 0.4920551 0.4866788 0.4860907 0.4912455 0.4940571

```

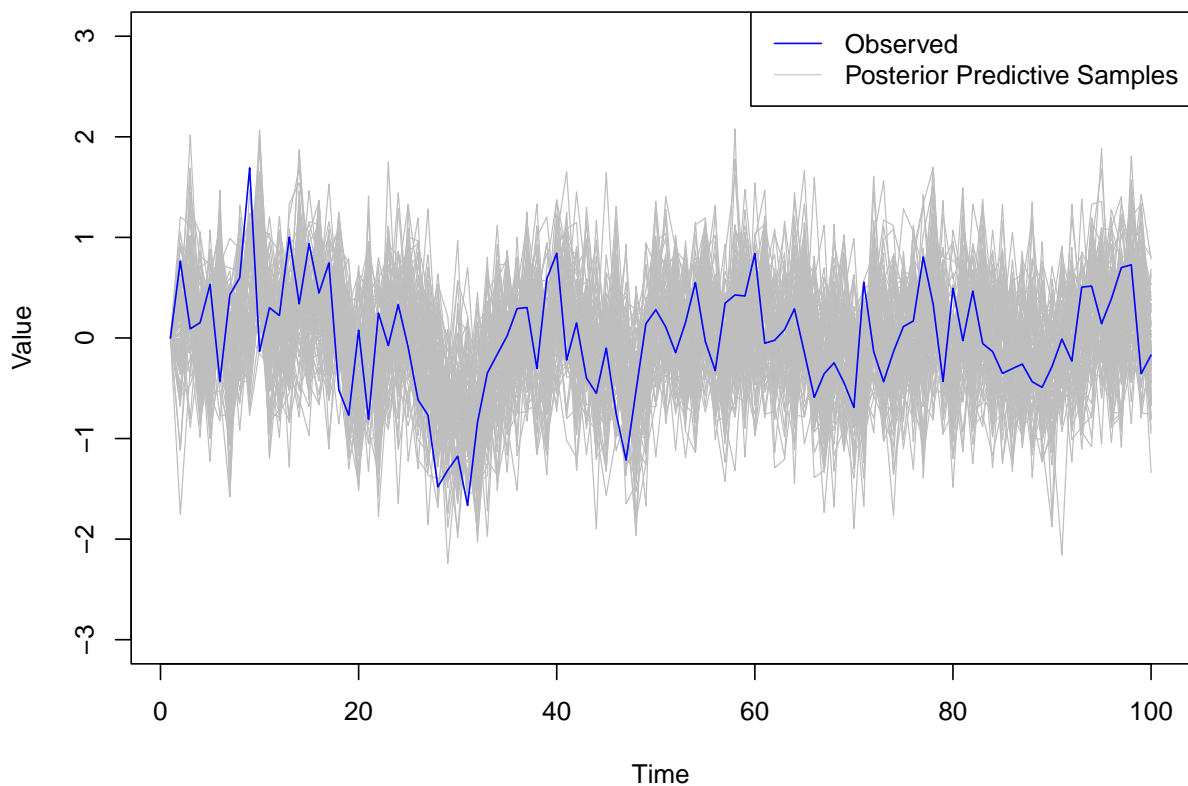
```
## [99] 0.4891187 0.4764067
```

7. (10 points) **Solution:**

The visualization of the posterior predictive samples against the observed data is as follows:

```
plot(
  x = seq_len(length(y)),
  y = samples[1, ],
  col = "gray",
  type = "l",
  lwd = 0.75,
  main = "Posterior Predictive Samples vs. Observed Data",
  xlab = "Time",
  ylab = "Value",
  ylim = c(-3.0, 3.0)
)
for (i in 2:100) {
  lines(x = seq_len(length(y)), y = samples[i, ], col = "gray", lwd = 0.75)
}
lines(x = seq_len(length(y)), y = y, col = "blue", lwd = 1.0)
legend(
  "topright",
  legend = c("Observed", "Posterior Predictive Samples"),
  col = c("blue", "gray"),
  lwd = c(1, 0.5),
  bg = "white",
)
```

**Posterior Predictive Samples vs. Observed Data**



From the visualization, we see that the posterior predictive samples have a wider range than the observed data, indicating that the model has a higher uncertainty in prediction and may not be able to capture the true data distribution well. Therefore, the model is **not a good fit** for the observed data. My expectation for a good model is that the posterior predictive samples are close to the observed data with a similar tendency and range.



## Real Data

1. (5 points) **Solution:**

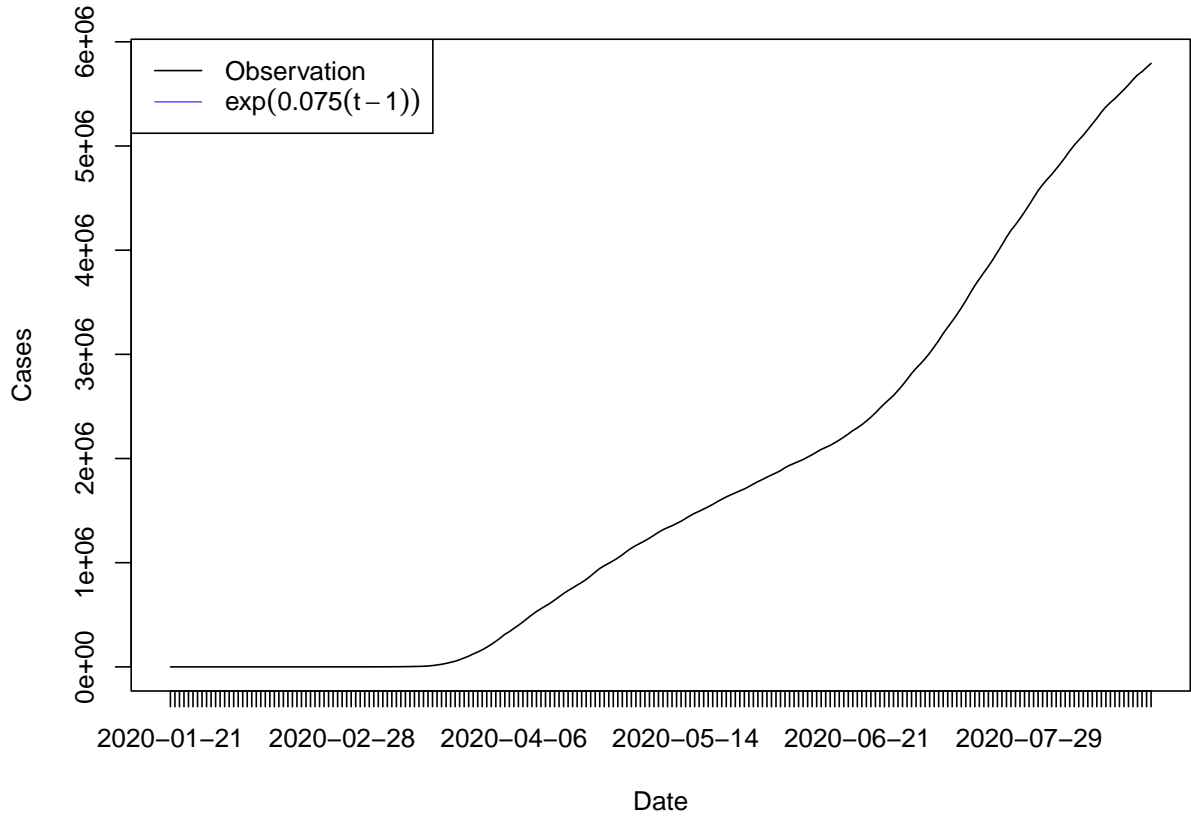
**No**, I do not believe that the information given is sufficient for analyzing the data. The GitHub page has only provided the dataset file, source of data, and definitions of different sorts of cases. All these pieces of information are related to the observation, i.e., COVID-19 cases, but none of them reflect any information to support the prior distribution of the parameters in the model. Other information such as demographic statistics, geofencing data, activity data, etc., can be useful for data analysis.

2. (5 points) **Solution:**

Before establishing the model, the raw data is visualized as follows:

```
# Read data
if (file.exists("data/covid_us.txt")) {
  data <- read.table("data/covid_us.txt", header = TRUE, sep = ",")
  y <- data[["cases"]]
} else {
  stop("File not found: 'computation_data_hw_1.csv' at ", getwd())
}
plot(
  x = seq_len(length(y)),
  y = y,
  type = "l",
  main = "COVID-19 Cases in the US",
  xlab = "Date",
  ylab = "Cases",
  xaxt = "n"
)
axis(side = 1, labels = data[["date"]], at = seq_len(length(y)))
legend(
  "topleft",
  legend = c("Observation", expression(y=exp(0.075 * (t - 1)))),
  col = c("black", "blue"),
  lwd = c(1, 0.5),
  bg = "white",
)
```

## COVID-19 Cases in the US



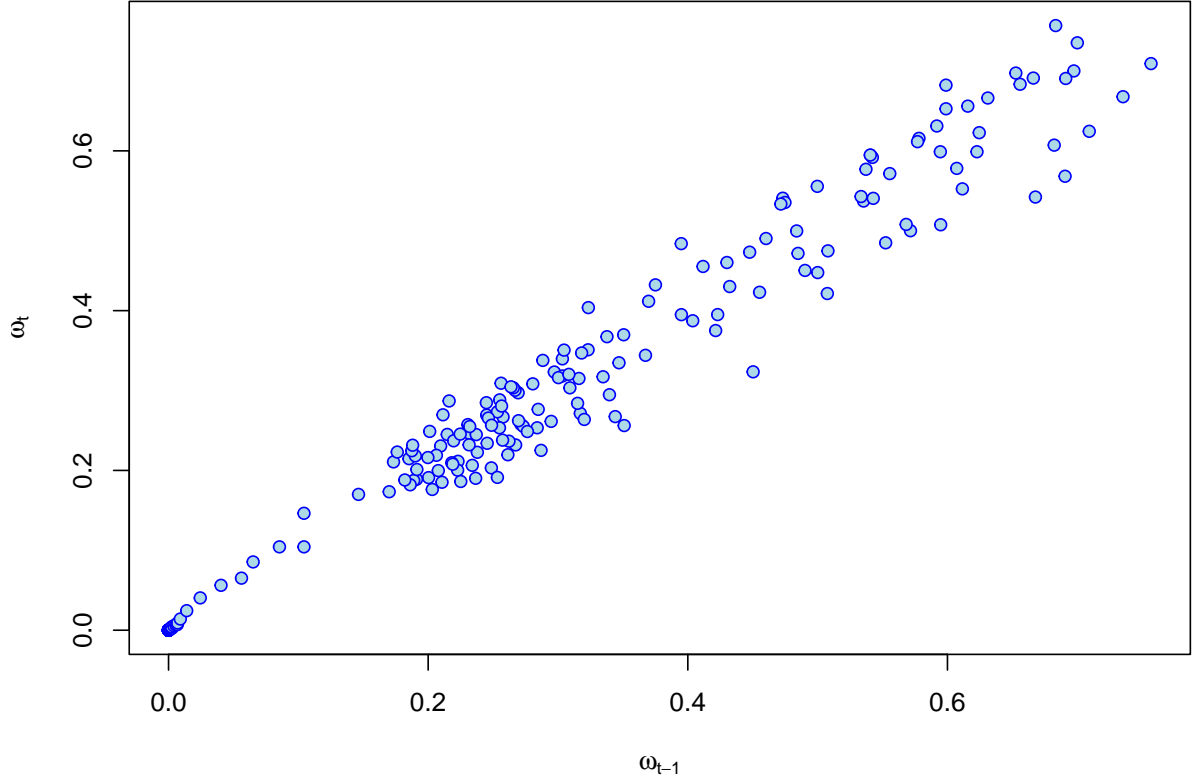
From the visualization, it is clear that the raw data has an exponential growth trend (blue line), which can not be directly captured by the AR(1) model. Therefore, the raw data needs to be transformed before being used for fitting an AR(1) model. The transformation consists of two steps: first, differentiating the transformed with `lag = 1` to filter out the non-stationary trend that is not captured by the AR(1) model, and then, scale the differentiated data by a factor of  $1/100000$  to make the data lying in a reasonable range:

$$\omega_t = \frac{1}{100000}(y_t - y_{t-1}).$$

For consistency,  $\omega_1 = 0$ . The visualization of the transformed data is as follows:

```
omg <- c(0, diff(y, lag = 1) / 100000)
plot(
  x = omg[seq(1, length(omg) - 1)],
  y = omg[seq(2, length(omg))],
  type = "p",
  pch = 21,
  col = "blue",
  bg = "lightblue",
  main = "Visualization of Transformed Data",
  xlab = expression(omega[t - 1]),
  ylab = expression(omega[t]),
)
```

## Visualization of Transformed Data



As shown in the visualization, the transformed data has a more linear relationship between consecutive observations. Finally, the AR(1) model is constructed as:

$$\omega_i = \rho\omega_{i-1} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2).$$

### 3. (5 points) **Solution:**

Based on the visualization above, a majority of the point cloud is growing linearly with a slope around 1.0 and slopes are non-negative. Therefore, the prior distribution of  $\rho$  can be set as a Gamma distribution  $\rho \sim \text{Gamma}(1, 1)$  with a mean of 1.0. For the log-variance  $\log(\sigma)$ , without losing generality, the prior can be set a normal distribution  $\log(\sigma) \sim \mathcal{N}(0, 1^2)$ .

### 4. (10 points) **Solution:**

With the above specified prior distributions, the log-posterior density is calculated as follows:

$$\begin{aligned} \log p(\rho, \log(\sigma) | \omega_1, \dots, \omega_n) &= \log L(\rho, \log(\sigma) | \omega_1, \dots, \omega_n) + \log p(\rho) + \log p(\log(\sigma)) + C \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega_i - \rho\omega_{i-1})^2 + \log \left( \frac{1}{\Gamma(1)} \exp(-\rho) \right) - \frac{\log(\sigma)^2}{2} + C' \\ &= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (\omega_i - \rho\omega_{i-1})^2 - \rho - \frac{\log(\sigma)^2}{2} + C', \end{aligned}$$

Therefore, the model fitted using data from the first time point until June 30 is as follows:

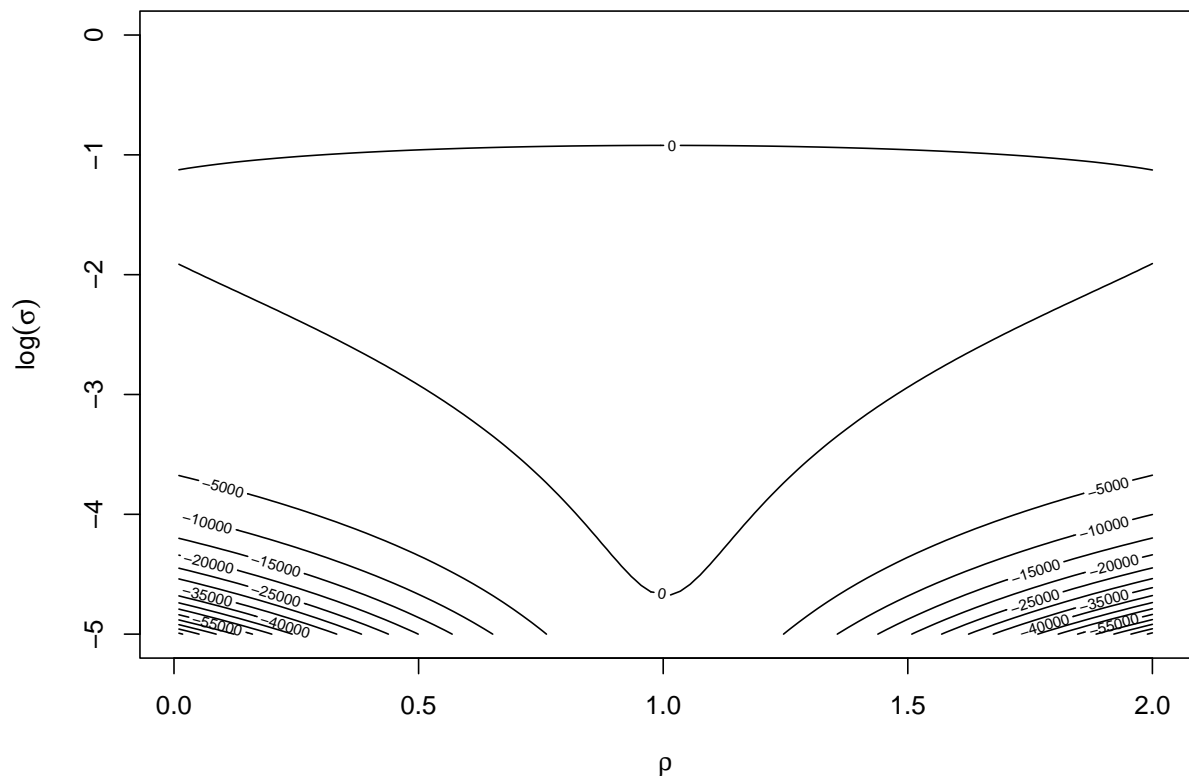
```
y <- c(
0,
diff(y[as.Date(data$date) < as.Date("2020-06-30")], lag = 1) / 100000
```

```

)
rho <- seq(0.01, 2.0, length = 100)
log_sig <- seq(-5.0, 0.0, length = 100)
log_lik <- outer(rho, log_sig, Vectorize(ar_loglik))
contour(
  x = rho,
  y = log_sig,
  z = log_lik,
  main = "Log-Likelihood Function",
  xlab = expression(rho),
  ylab = expression(log(sigma)),
  nlevels = 20,
)

```

**Log-Likelihood Function**



The visualization of the logarithm of unnormalized posterior density is shown below:

```

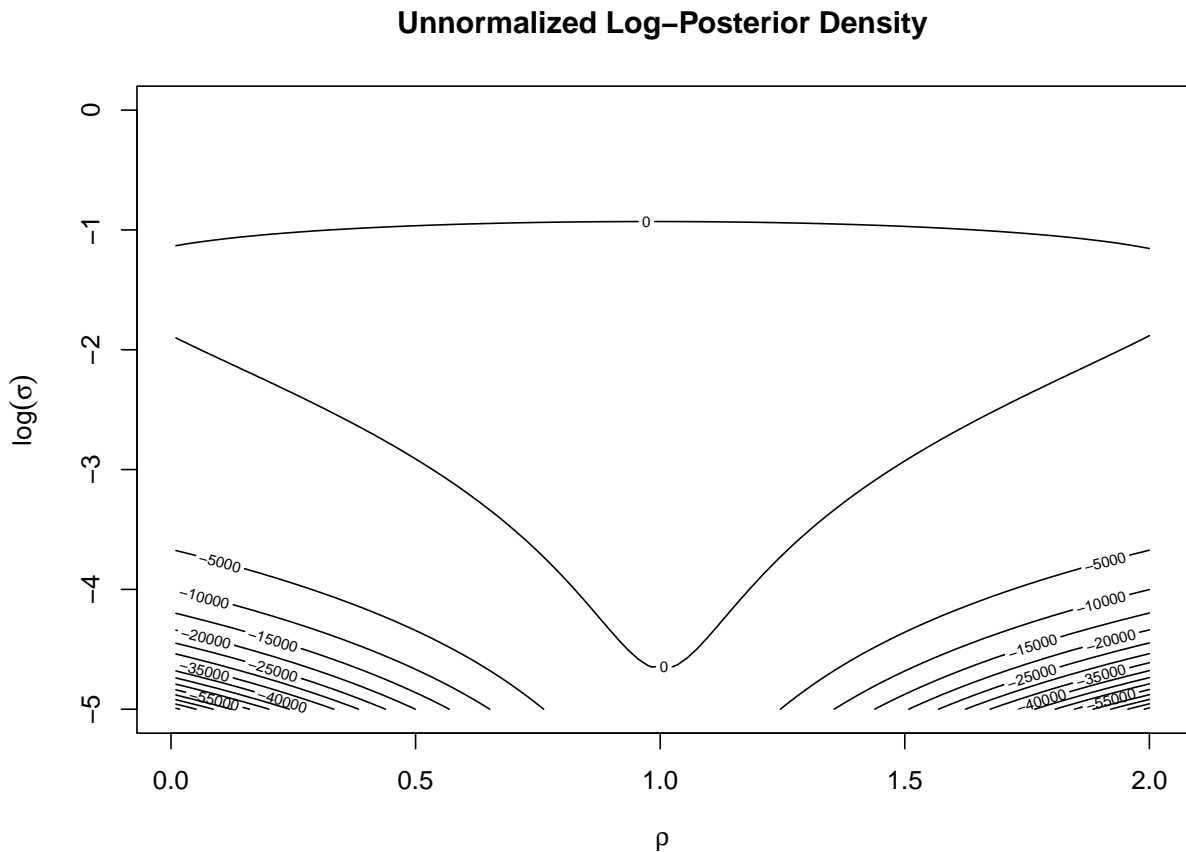
ar_logpost <- function(rho, log_sig) {
  log_prior <- -rho - 0.5 * (log_sig * log_sig)
  log_post <- ar_loglik(rho, log_sig) + log_prior
  return(log_post)
}
log_post <- outer(rho, log_sig, Vectorize(ar_logpost))
contour(
  x = rho,
  y = log_sig,
  z = log_post,
  main = "Unnormalized Log-Posterior Density",

```

```

xlab = expression(rho),
ylab = expression(log(sigma)),
nlevels = 20,
)

```



Similarly, I draw 1000 samples from the grid and calculate the summary statistics as follows:

```

library(moments)

grid_log_post <- outer(rho, log_sig, Vectorize(ar_logpost))
# Calculate the normalized probability density
probs <- exp(grid_log_post - max(grid_log_post))
probs <- probs / sum(probs)
print(dim(probs))

## [1] 100 100

# Randomly sample from the grid
indices <- sample(
  x = seq_len(length(as.vector(probs))),
  size = 1000,
  replace = TRUE,
  prob = probs
)
rho_sample <- rho[((indices - 1) %% nrow(probs)) + 1]
log_sig_sample <- log_sig[((indices - 1) %% nrow(probs)) + 1]

# Calculate summaries for rho

```

```

print(quantile(rho_sample, probs = c(0.025, 0.25, 0.5, 0.75, 0.975)))

##      2.5%      25%      50%      75%      97.5%
## 0.9949495 0.9949495 0.9949495 1.0150505 1.0150505

sprintf("Mean: %.4f", mean(rho_sample))

## [1] "Mean: 1.0031"

sprintf("Standard Deviation: %.4f", sd(rho_sample))

## [1] "Standard Deviation: 0.0105"

sprintf("Skewnewss: %.4f", skewness(rho_sample))

## [1] "Skewnewss: 0.3297"

sprintf("Kurtosis: %.4f", kurtosis(rho_sample))

## [1] "Kurtosis: 1.9140"

# Calculate summaries for log(sigma)
print(quantile(log_sig_sample, probs = c(0.025, 0.25, 0.5, 0.75, 0.975)))

##      2.5%      25%      50%      75%      97.5%
## -3.737374 -3.686869 -3.636364 -3.636364 -3.535354

sprintf("Mean: %.4f", mean(log_sig_sample))

## [1] "Mean: -3.6495"

sprintf("Standard Deviation: %.4f", sd(log_sig_sample))

## [1] "Standard Deviation: 0.0567"

sprintf("Skewness: %.4f", skewness(log_sig_sample))

## [1] "Skewness: 0.0189"

sprintf("Kurtosis: %.4f", kurtosis(log_sig_sample))

## [1] "Kurtosis: 2.9988"

```

5. (10 points) **Solution:**

Similarly, the posterior predictive samples are compared to the observed transformed data. The code implementation is as follows:

```

params <- data.frame(rho = rho_sample, log_sig = log_sig_sample)
samples <- matrix(NA, nrow = nrow(params), ncol = length(y))
for (i in seq_len(nrow(params))) {
  samples[i, ] <- ar_post_predictive(params[i, "rho"], params[i, "log_sig"])
}

# Visualize scatter plot
plot(
  x = seq_along(samples[1, seq(1, length(y))]),
  y = samples[1, seq(1, length(y))],
  type = "l",
  col = "gray",
  xlab = "Time",
  ylab = expression(omega[t]),

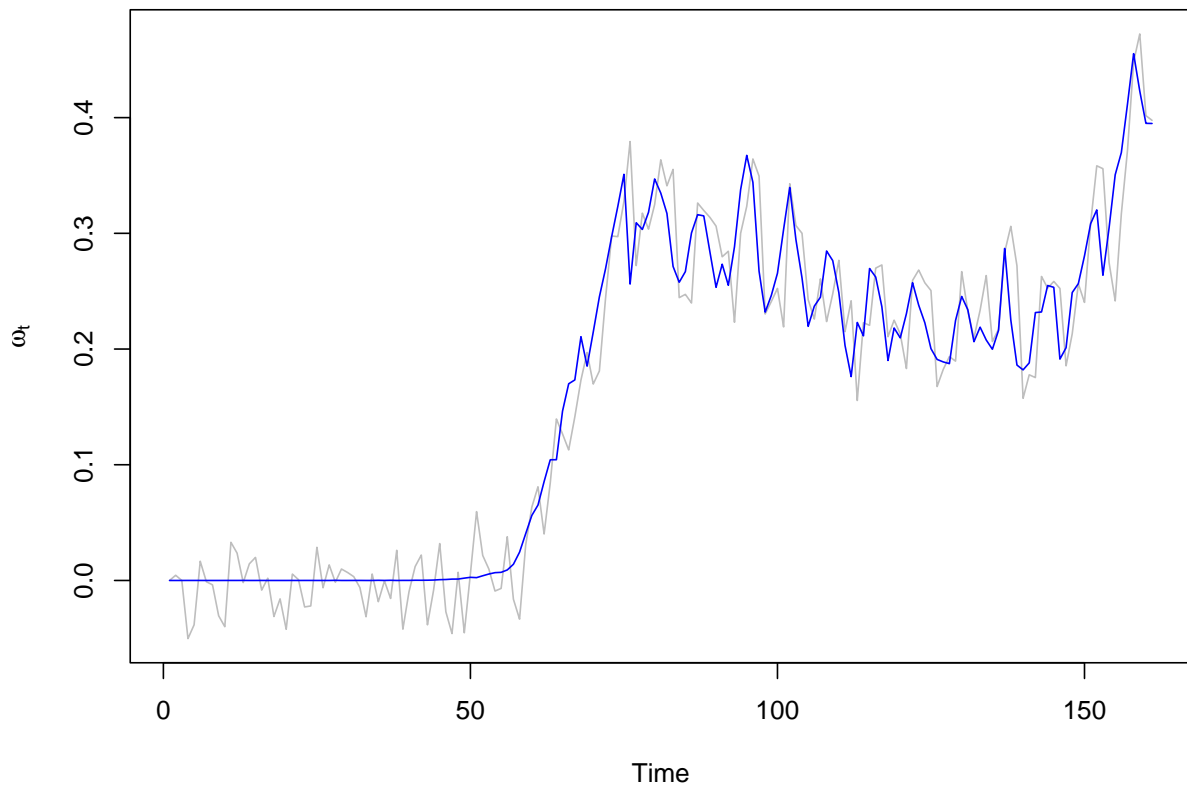
```

```

    main = "Posterior Predictive Samples vs. Observed Data",
  )
  for (i in 2:1000) {
    lines(x = seq_along(samples[i]), y = samples[i], col = "gray")
  }
  lines(x = seq_along(y), y = y, col = "blue")

```

### Posterior Predictive Samples vs. Observed Data



Following the posterior predictive check covered in the lecture, the visualization of the skewness and kurtosis of the posterior predictive samples against the observed data is as follows:

```

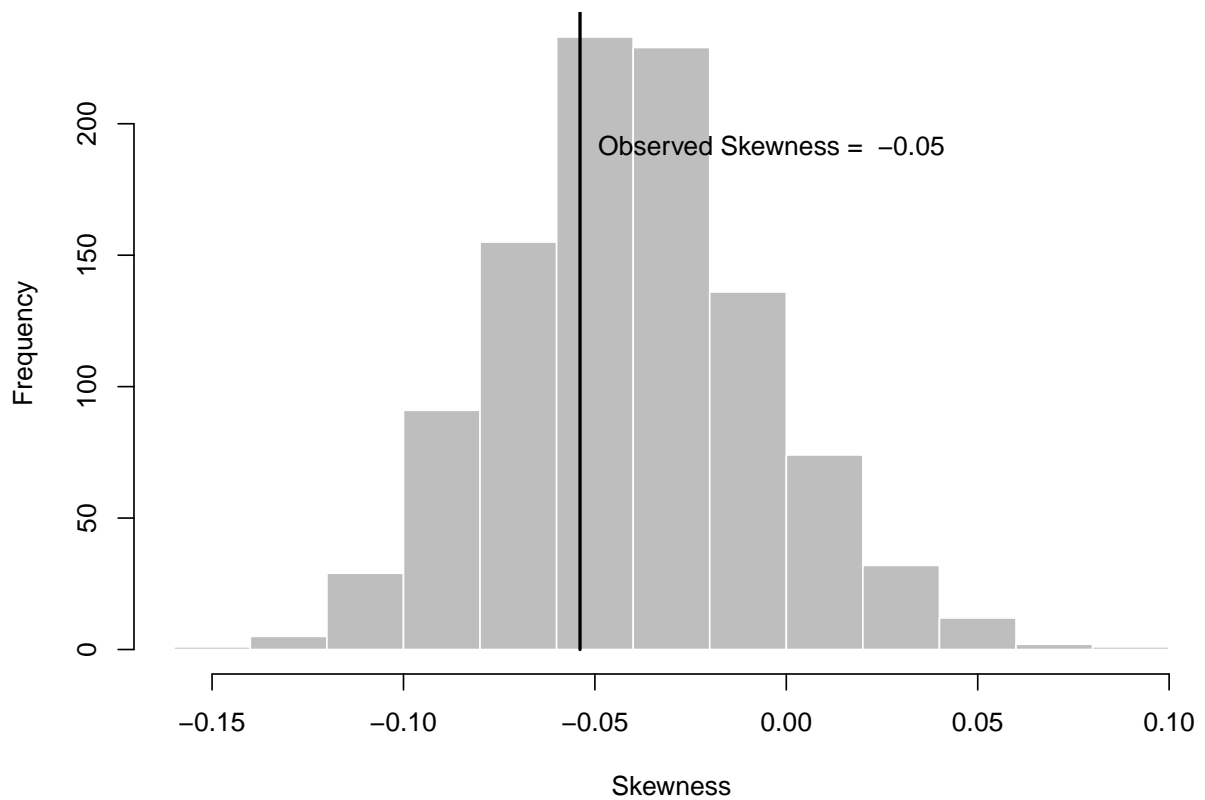
skewness_obs <- skewness(y[seq(2, length(y))])
sample_skewness <- apply(samples, 1, function(x) skewness(x[seq(2, length(y))]))
kurtosis_obs <- kurtosis(y[seq(2, length(y))])
sample_kurtosis <- apply(samples, 1, function(x) kurtosis(x[seq(2, length(y))]))

hist(
  x = sample_skewness,
  main = "Skewness of Posterior Predictive Samples",
  xlab = "Skewness",
  col = "gray",
  border = "white",
)
lines(
  x = c(skewness_obs, skewness_obs),
  y = c(0, 300),
  col = "black",
  lwd = 2,

```

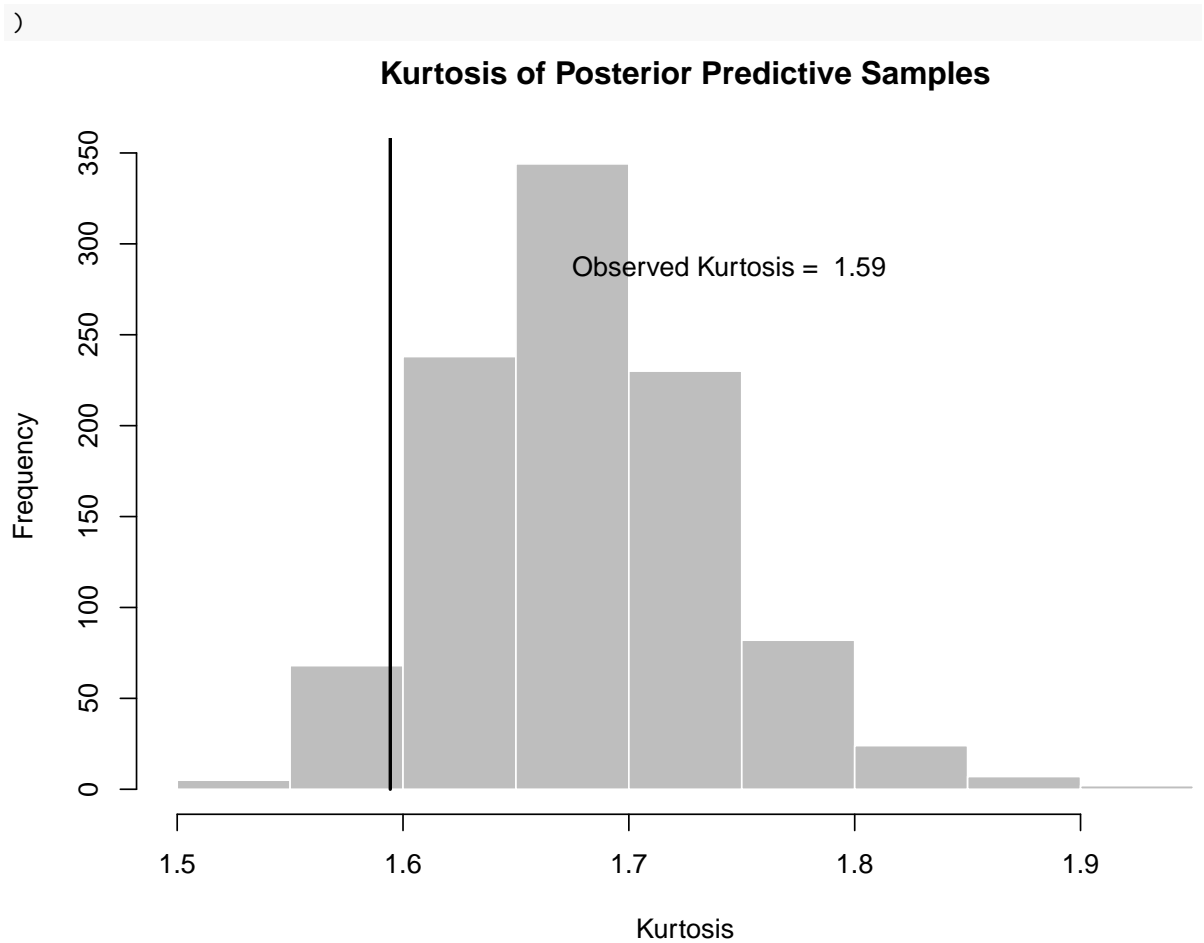
```
)
text(
  x = skewness_obs + 0.05,
  y = 200,
  labels = paste("Observed Skewness = ", round(skewness_obs, 2)),
  pos = 1
)
```

### Skewness of Posterior Predictive Samples



```
hist(
  x = sample_kurtosis,
  main = "Kurtosis of Posterior Predictive Samples",
  xlab = "Kurtosis",
  col = "gray",
  border = "white",
)
lines(
  x = c(kurtosis_obs, kurtosis_obs),
  y = c(0, 500),
  col = "black",
  lwd = 2,
)
text(
  x = kurtosis_obs + 0.15,
  y = 300,
  labels = paste("Observed Kurtosis = ", round(kurtosis_obs, 2)),
  pos = 1
)
```





Based on the posterior predictive checks, the model provides good predictions for the data until June 30. The skewness and kurtosis of the observed data are close to the expected values from the posterior predictive samples, and the time-series visualization shows that the posterior predictive samples are close to the observed data.

6. (10 points) Simulate posterior predictions for the cases and deaths in the United States from July 1 until August 25. Compare these posterior predictions to the actual values of the two variables, and comment on the quality of the AR(1) model for predicting future cases and deaths in this context.
7. (10 points) Now replicate the analyses you performed in step 4 on the entire United States for the state of Indiana individually. Discuss if you think the results agree with the results of the entire United States.
8. (5 points) Provide a non-technical explanation of your finds for an audience with minimal statistical training. Specifically, describe the AR(1) model in lay terms, explain whether or not this model provides good predictions in this context, and provide intuitive justifications for why the AR(1) model succeeds or fails in this context.