

R handout: Monte Carlo Method and Application

Yuhui Yao*

This handout will show various ways to conduct the Likelihood ratio test in practice. Suppose we are interested in the following hypothesis testing for the mean of sepal width of Virginica (from the Iris data set (Fisher (1936))) $H_0 : \mu = \mu_0 = 3$ vs. $H_1 : \mu = \mu_1 \neq 3$ under the assumption that the observations X_i 's are i.i.d. the normal distribution with $N(\mu, \sigma^2 = 0.1040)$. Also, let the significance level be $\alpha = 0.05$. First of all, I am going to describe the sepal width of Virginica.

- R code

```
# set a seed
set.seed(12345)

# specify the address of the data
data.addr <- 'https://raw.githubusercontent.com/
  bolus123/R-handout/master/LikelihoodRatioTest/iris.csv'

# load the data
data <- read.csv(file = data.addr)

# name the columns
names(data) <- c(
  'sl' #sepal length
  , 'sw' #sepal width
  , 'pl' #petal length
  , 'pw' #petal width
  , 'class' #class
)

# get the fraction of Virginica
data.virginica <- data[data$class == 'Iris-virginica', ]

# describe the Virginica data
mean(data.virginica$sw)

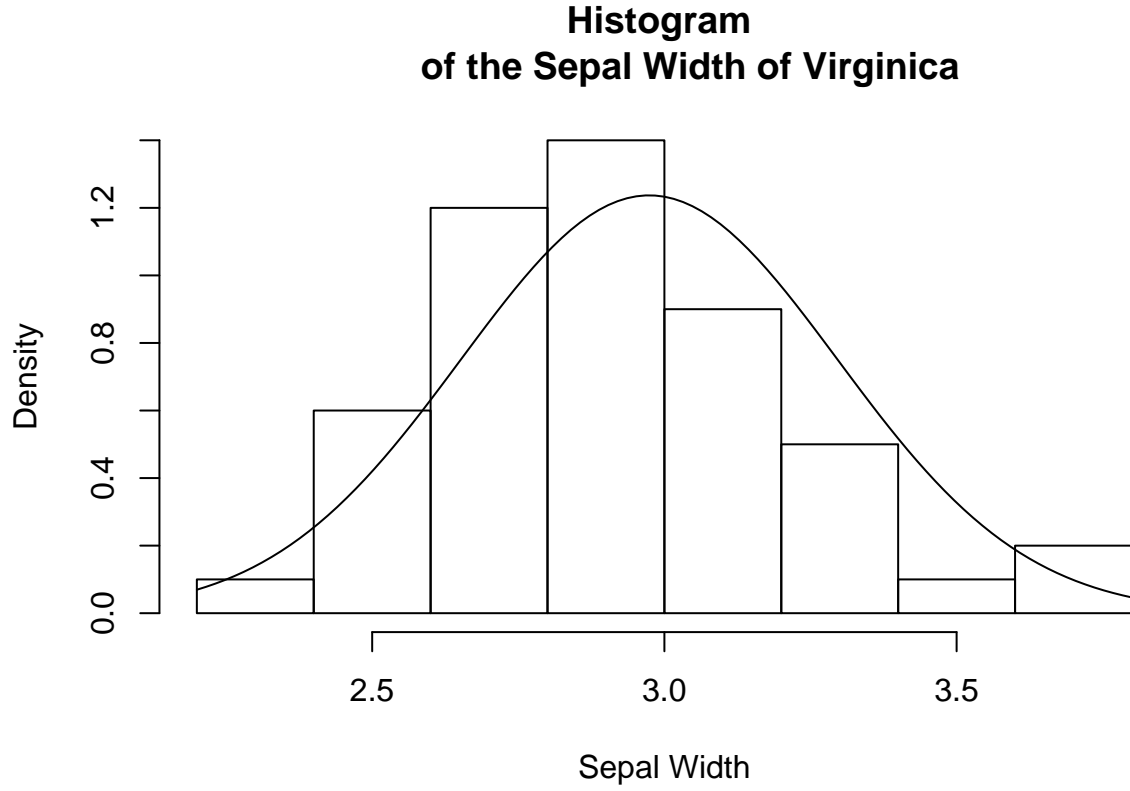
## [1] 2.974

var(data.virginica$sw)

## [1] 0.1040041

# build a histogram
hist(data.virginica$sw, freq = FALSE, main = 'Histogram
  of the Sepal Width of Virginica', xlab = 'Sepal Width')
# add a fitted line
curve(dnorm(x, mean(data.virginica$sw), sd(data.virginica$sw)), add = TRUE)
```

*The University of Alabama



1. Exact Method

The likelihood

$$L(\mu; \underline{x}) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} e^{-(x_i - \mu)^2 / (2\sigma^2)}$$

The log-likelihood

$$\ln L(\mu; \underline{x}) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \quad (1)$$

It is easy to show the m.l.e. be $\hat{\mu} = \bar{x}$. The log-likelihood ratio

$$\begin{aligned} \Lambda &= \ln\left(\frac{L(\mu_0; \underline{x})}{L(\hat{\mu}; \underline{x})}\right) = \ln L(\mu_0; \underline{x}) - \ln L(\hat{\mu}; \underline{x}) \\ &= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu_0)^2 + \left(\frac{n}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \hat{\mu})^2\right) \end{aligned}$$

$$\begin{aligned}
\sum_{i=1}^n (x_i - \mu_0)^2 &= \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu_0)^2 \\
&= \sum_{i=1}^n [(x_i - \bar{x})^2 + 2(x_i - \bar{x})(\bar{x} - \mu_0) + (\bar{x} - \mu_0)^2] \\
&= \sum_{i=1}^n [(x_i - \bar{x})^2] + n(\bar{x} - \mu_0)^2
\end{aligned}$$

So

$$\Lambda = -\frac{n}{2\sigma^2}(\bar{x} - \mu_0)^2$$

Because $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$, $\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1) \Rightarrow (\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}})^2 \sim \chi_1^2$, according to the equal-tail assumption, the size

$$P(-2\Lambda \leq C_1 \cup -2\Lambda \geq C_2 | H_0) = P(-2\Lambda \leq C_1 \cup -2\Lambda \geq C_2 | \mu = \mu_0) = \alpha$$

where under the null hypothesis $-2\Lambda \sim \chi_1^2$. Wilks (1938) proved that, if the sample size is large enough and H_0 is true, for a nested model, -2Λ is asymptotically following the chi-squared distribution with the degrees of freedom approximately equal to the difference between the degrees of freedom under the alternative hypothesis and the degrees of freedom under the null hypothesis. This result is called **Wilks' Theorem**. Due to our -2Λ exactly following the chi-squared distribution with 1 degree of freedom, I am not going to show the approximation which is as same as the exact method. If you are interested in the approximation, you can try to assume the sepal width follows a gamma distribution, and then you may need to use the approximation.

- R code

```

# sample size of Virginica
n <- dim(data.virginica)[1]

# sigma2 is known
sigma2 <- 0.1040
sigma <- sqrt(sigma2)

# under the null hypothesis
mu.0 <- 3

# under the alternative hypothesis
mu.1 <- mean(data.virginica$sw)

# calculate the statistic
Lambda <- - n / 2 / sigma2 * (mu.1 - mu.0)^2
Lambda <- -2 * Lambda
Lambda

## [1] 0.325

# critical values for alpha = 0.05 under the equal-tailed assumption
c1 <- qchisq(0.025, 1) # critical value for the lower tail
c1

## [1] 0.0009820691

```

```

c2 <- qchisq(0.975, 1) # critical value for the upper tail
c2

## [1] 5.023886
# calculate p-value
p <- pchisq(Lambda, 1)
p.value <- ifelse(p > 0.5, 1 - (1 - p) * 2, p * 2 )
p.value

## [1] 0.8627636

```

2. Bootstrap Method

Suppose the analytical form of -2Λ is too complicated and we do not have enough sample size to reach the asymptote but we can calculate the likelihoods. We can do our test by simulating the empirical distribution of the log-likelihood ratio -2Λ .

(a) Parametric Bootstrap

We will **simulate the data for a specific distribution under the null hypothesis** and then do the test repeatedly.

Steps for our problem:

- i. **Simulate a sample from the given distribution under the null hypothesis**
- ii. Calculate the likelihood under the null hypothesis
- iii. Calculate the likelihood under the alternative hypothesis (Notice that the m.l.e. need to be re-estimated)
- iv. Calculate and save $-2\Lambda_s$ for the empirical distribution of -2Λ
- v. repeat (1) - (4) until the process reaches the maximum of simulations
- vi. Compare the -2Λ calculated by the original data with the empirical distribution of -2Λ

- R code

```

# maximum number of simulations
sim <- 10000

# sample size of Virginica
n <- dim(data.virginica)[1]

# get the empirical distribution based on the parametric bootstrap
ref.par <- rep(NA, sim)

for (i in 1:sim){

  # simulate data under the null hypothesis
  X <- rnorm(n, mu.0, sigma)

```

```

# calculate Lambda
lnL.0 <- sum(log(dnorm(X, mu.0, sigma)))
lnL.1 <- sum(log(dnorm(X, mean(X), sigma)))
ref.par[i] <- -2 * (lnL.0 - lnL.1)

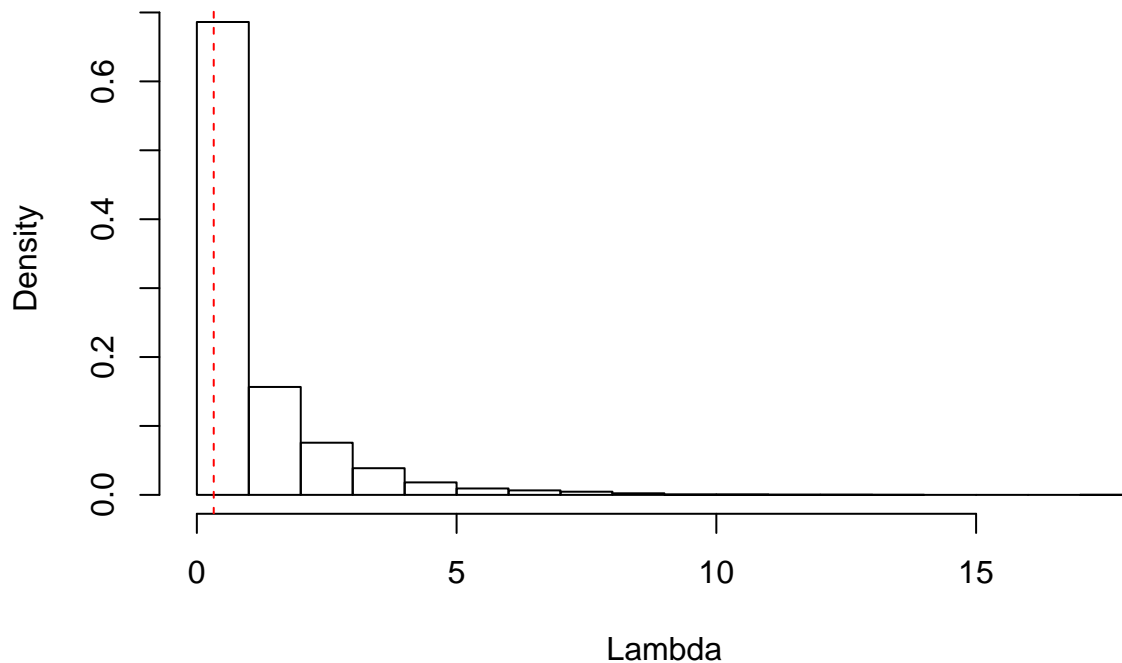
}

# describe the empirical distribution of Lambda based on the parametric bootstrap
hist(ref.par, freq = FALSE, main = 'Histogram of the empirical distribution of Lambda\n
    based on the parametric bootstrap', xlab = 'Lambda')
# point out the Lambda from the original data
abline(v = Lambda, lty = 2, col = 'red')

```

Histogram of the empirical distribution of Lambda

based on the parametric bootstrap



```

# critical values for alpha = 0.05 under the equal-tailed assumption
c1 <- quantile(ref.par, 0.025) # critical value for the lower tail
c1

##          2.5%
## 0.001088818

c2 <- quantile(ref.par, 0.975) # critical value for the upper tail
c2

##          97.5%

```

```
## 4.989207
# calculate p-value
p <- mean(ref.par < Lambda)
p.value.par <- ifelse(p > 0.5, 1 - (1 - p) * 2, p * 2)
p.value.par

## [1] 0.8748
```

(b) Nonparametric Bootstrap

We will **resample the data from the original data set** and then do the test repeatedly.

Steps for our problem:

- i. **Resample a sample from the original data**
- ii. the rest of steps are exactly same as the ones ii. - vi. for the parametric bootstrap shown above.

```
# maximum number of simulations
sim <- 10000

# sample size of Virginica
n <- dim(data.virginica)[1]

# get the empirical distribution based on the nonparametric bootstrap
ref.nonpar <- rep(NA, sim)

for (i in 1:sim){

  # resample from the original data set
  X <- sample(data.virginica$sw, size = n, replace = TRUE)

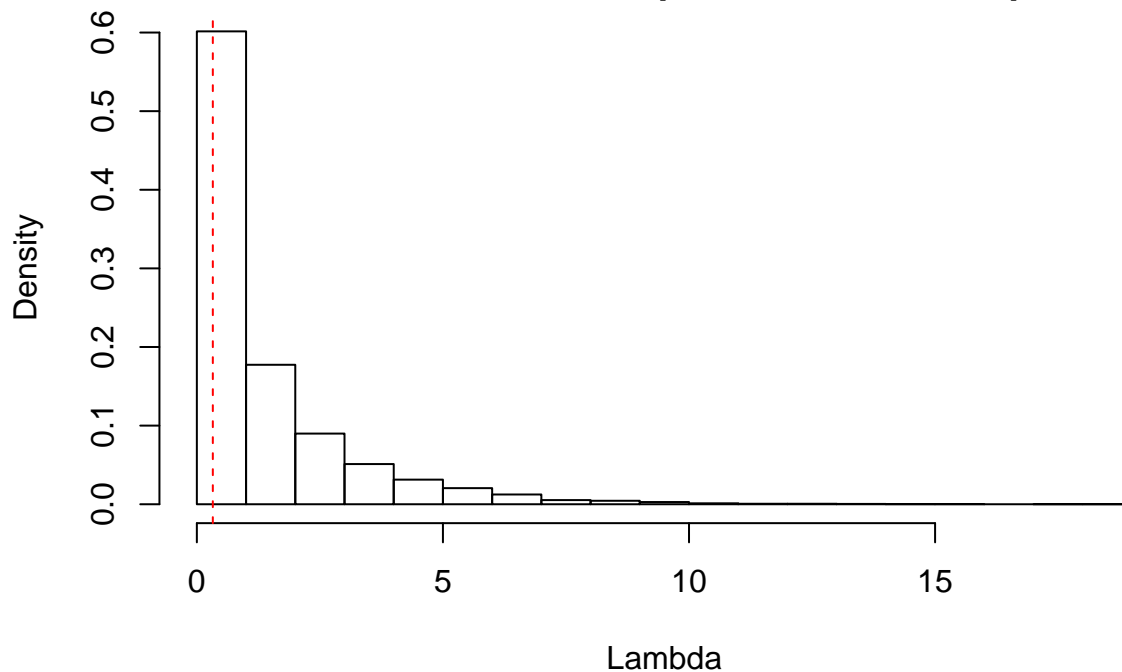
  # calculate Lambda
  lnL.0 <- sum(log(dnorm(X, mu.0, sigma)))
  lnL.1 <- sum(log(dnorm(X, mean(X), sigma)))
  ref.nonpar[i] <- -2 * (lnL.0 - lnL.1)

}

# describe the empirical distribution of Lambda based on the nonparametric bootstrap
hist(ref.nonpar, freq = FALSE, main = 'Histogram of the empirical distribution of Lambda\n
    based on the nonparametric bootstrap', xlab = 'Lambda')
# point out the Lambda from the original data
abline(v = Lambda, lty = 2, col = 'red')
```

Histogram of the empirical distribution of Lambda

based on the nonparametric bootstrap



```
# critical values for alpha = 0.05 under the equal-tailed assumption
c1 <- quantile(ref.nonpar, 0.025) # critical value for the lower tail
c1
```

```
##          2.5%
## 0.001923077
```

```
c2 <- quantile(ref.nonpar, 0.975) # critical value for the upper tail
c2
```

```
##          97.5%
## 6.248077
```

```
# calculate p-value
p <- mean(ref.nonpar < Lambda)
p.value.nonpar <- ifelse(p > 0.5, 1 - (1 - p) * 2, p * 2)
p.value.nonpar
```

```
## [1] 0.7198
```

3. Comparison

The following is the table for the p-values from different methods

```
##          ExactMethod ParametricBootstrap NonparametricBootstrap
```

| | | | | |
|----|---------|--------|--------|--------|
| ## | p-value | 0.8628 | 0.8748 | 0.7198 |
|----|---------|--------|--------|--------|

The p-values from the exact method are similar with the the one from the Parametric Bootstrap method which is greater than the one from the Nonparametric Bootstrap method. So why do we have this result?

4. Practice (Goodness-of-fit)

Suppose we are interested in the following hypothesis testing for deciding whether the sample of sepal width of Virginica is from the specific normal distribution with mean $\mu = 3$ and variance $\sigma^2 = 0.1$. Follow the three methods of conducting the likelihood ratio test above, what is your setting for the test and what is your conclusion? (**Hint:** the degrees of freedom for this test is not equal to 1. Also, you need to estimate the m.l.e. of sample mean and sample variance, simultaneously).

Reference

Fisher, Ronald A. 1936. "The Use of Multiple Measurements in Taxonomic Problems." *Annals of Human Genetics* 7 (2). Wiley Online Library:179–88.

Wilks, Samuel S. 1938. "The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses." *The Annals of Mathematical Statistics* 9 (1). JSTOR:60–62.