

HW3

Hannah Zmuda

11/18/2020

Question 1

Problem

Use Monte Carlo simulation to investigate whether the empirical Type I error rate of the t-test is approximately equal to the nominal significance level α , when the sampled population is non-normal. The t-test is robust to mild departures from normality. Discuss the simulation results for the cases where the sampled population is (i) $\chi^2(1)$, (ii) Uniform(0,2), and (iii) Exponential(rate=1). In each case, test $H_0 : \mu = \mu_0$ vs $H_0 : \mu \neq \mu_0$, where μ_0 is the mean of $\chi^2(1)$, Uniform(0,2), and Exponential(1), respectively.

Solution

```
#Part a: looking at a type I error
#all follow a similar algorithm as shown on page 193 in SCRR
n <- 100 #number of replicates
a <- 0.05 #significance level alpha
muA <- mean(rchisq(n, df = 1)) #mu0 in part a
muB <- mean(runif(n, 0, 2)) #mu0 in part b
muC <- mean(rexp(n, rate = 1)) #mu0 in part c
#become alternatives in the estimate power of a test (b).

m <- 1000 #number of replicates
pA <- numeric(m)
pB <- numeric(m)
pC <- numeric(m)

for(i in 1:m){
  xA <- rchisq(n, df = 1) #sample dist part a
  xB <- runif(n, 0, 2) #sample dist part b
  xC <- rexp(n, rate = 1) #sample dist part c

  ttestA <- t.test(xA, alternative = "two.sided", mu = muA)
  ttestB <- t.test(xB, alternative = "two.sided", mu = muB)
  ttestC <- t.test(xC, alternative = "two.sided", mu = muC)

  pA[i] <- ttestA$p.value
  pB[i] <- ttestB$p.value
  pC[i] <- ttestC$p.value
}

pHatA <- mean(pA <= a)
pHatB <- mean(pB <= a)
pHatC <- mean(pC <= a)
```

```

seHatA <- sqrt(pHatA * (1- pHatA)/m)
seHatB <- sqrt(pHatB * (1- pHatB)/m)
seHatC <- sqrt(pHatC * (1- pHatC)/m)
#Means:
cat("The means (of the p-values) are: ", c(pHatA, pHatB, pHatC), "\n")

## The means (of the p-values) are:  0.288 0.125 0.545

#SE:
cat("The standard errors (of the p-values) are: ", c(seHatA, seHatB, seHatC))

## The standard errors (of the p-values) are:  0.01431978 0.01045825 0.01574722

#Part b: Estimating power of a test and outputting empirical power curves of the t-test from the three ;
#initial data array
mu1 <- c(seq(0,2,1/20))
MC <- length(mu1)
powA <- numeric(MC)
powB <- numeric(MC)
powC <- numeric(MC)
powNorm <- numeric(MC)

#select theta_1 from the parameter subspace
#muA, muB, and muC values
#set for loop
for(j in 1:MC){
  #Chi Squared Distribution
  pvalA <- replicate(MC, expr = {
    xA <- rchisq(n, df = 1) #sample dist part a
    ttestA <- t.test(xA, alternative = "two.sided",mu = mu1[j])
    ttestA$p.value
  })
  powA[j] <- mean(pvalA <= a)
  #Uniform Distribution
  pvalB <- replicate(MC, expr = {
    xB <- runif(n, 0, 2) #sample dist part b
    ttestB <- t.test(xB, alternative = "two.sided",mu = mu1[j])
    ttestB$p.value
  })
  powB[j] <- mean(pvalB <= a)
  #Exponential Distribution
  pvalC <- replicate(MC, expr = {
    xC <- rexp(n, rate = 1) #sample dist part c
    ttestC <- t.test(xC, alternative = "two.sided",mu = mu1[j])
    ttestC$p.value
  })
  powC[j] <- mean(pvalC <= a)
  #Normal Distribution using power.t.test function
  powerTest <- power.t.test(n = MC, delta = mu1[j],sig.level = a,alternative = "two.sided")
  powNorm[j] <- powerTest$power
}

#making of the data frame
mean <- c(mu1,mu1,mu1,mu1)

```

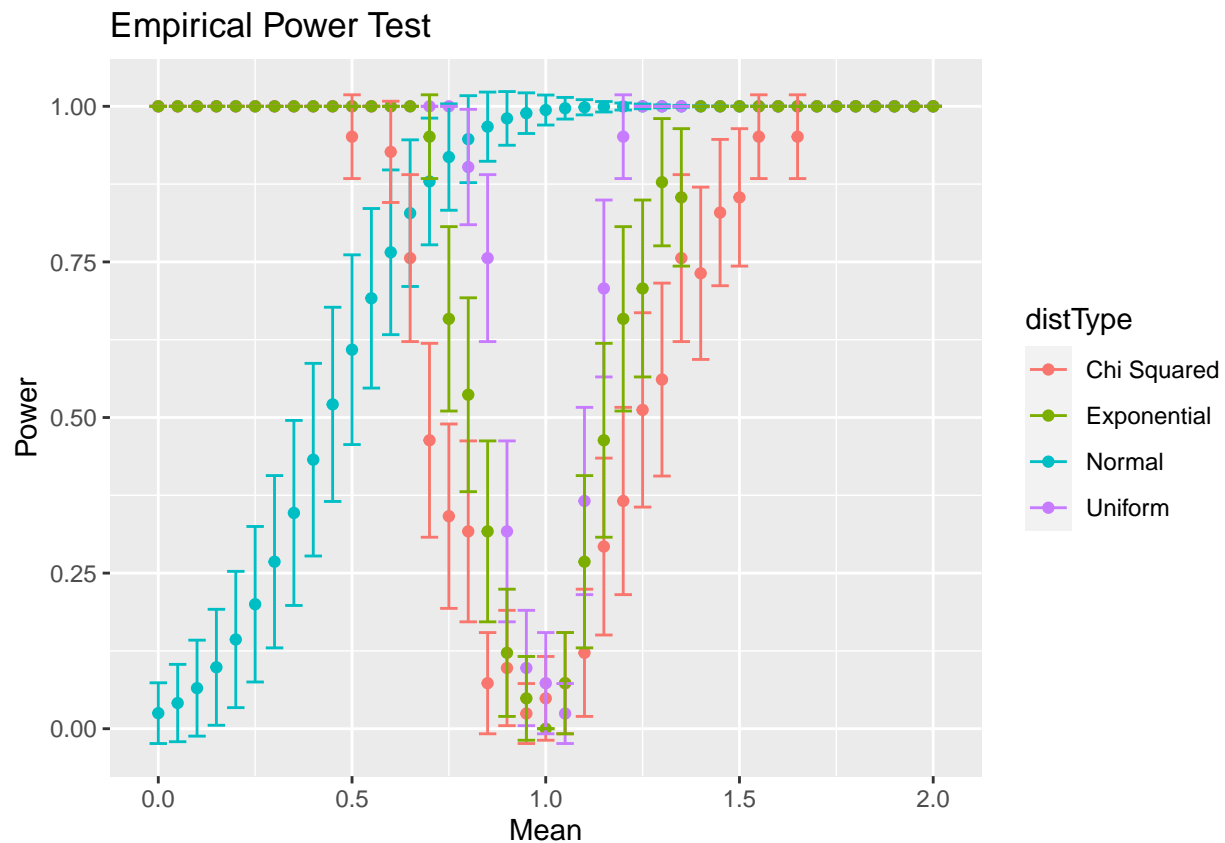
```

powerR <- c(powNorm, powA, powB, powC)
distType <- c(replicate(MC, "Normal"),replicate(MC, "Chi Squared"),replicate(MC, "Uniform"),replicate(MC, "Exponential"))

data <- data.frame(mean, powerR,distType)

ggplot(data = data, aes(x = mean, y = powerR, color = distType)) +
  geom_point() +
  labs(x = 'Mean', y = 'Power', title = 'Empirical Power Test') +
  geom_errorbar(data = data, mapping = aes(x = mean, ymin = powerR - 2*(sqrt(powerR * (1 - powerR)/MC)),

```



Question 2

Part a and b: MC and IS

For part (a) and (b), we are comparing Monte Carlo method and Importance Sampling (IS) to estimate α using the Z-test. In order to do Importance Sampling in part (b), we must first derive the weights: $f(X)/g(X)$. The problem statement tells us that the function $g(x) = \text{Pois}(1.5\lambda)$, which means the importance function, $f(x)$, is also based on the Poisson pdf. Therefore,

$$f(x) = \frac{e^\lambda \lambda^x}{x!}$$

and

$$g(x) = \frac{e^{-1.5\lambda} 1.5\lambda^x}{x!}$$

```
set.seed(475)
#part a estimate alpha using MCEM
n <- 100 # distribution size
m <- 100 #Monte Carlo sample size
l <- 2 #lambda is equal to 2
#z test with Monte Carlo
ztest <- replicate(m, expr = {
  x <- rpois(n, l) #samples
  ztest <- (mean(x)-2)/(sd(x)/sqrt(n))
})
#alpha estimate
aMC <- mean(ztest > 2.326)
cat("Alpha estimate using MC:", aMC, "\n")

## Alpha estimate using MC: 0.02

#part b: estimate alpha using Importance Sampling
g <- function(x) (exp(-1.5*l)*((1.5*l)^x))/factorial(x)
f <- function(x) (exp(-1)*(l^x))/factorial(x)
phi <- function(x) as.numeric((mean(x)-2)/(sd(x)/sqrt(n)) > 2.326)
weight <- function(x) f(x) / g(x)
out <- numeric(m)

out <- replicate(m, expr = {
  x <- rpois(n, l) #samples
  out <- phi(x)*weight(x)
})

isOut <- mean(out)
cat("Alpha estimate using IS:", isOut, "\n")

## Alpha estimate using IS: 0.02424404
```

Part c: Which is better?

Both outputs of the function are off from the expected 0.05 estimate of α . I would argue that the Monte Carlo is easier to implement because the weights do not need to be derived. However, IS estimates the α value closer to the expected value, compared to using Monte Carlo by itself.

Question 3

```
salmon <-read.table("salmon.dat", header = TRUE)
```

Question 4

Question 5