

# Statistical Inference Class Project

Cathy Snell

October 10, 2018

## Part 1: Simulation Exercise

### Overview

We will investigate the exponential distribution in R and compare it with the Central Limit Theorem (CLT). The CLT states that “the distribution of averages of iid random variables becomes that of a standard normal as the sample size increases.” (from 07 02 Asymptotics and the CLT lecture transcript). We will show that the mean of the sample means approaches the theoretical mean of the exponentials, is distributed as a normal curve, and that the variance decreases as the sample size increases.

### Simulations

We have been directed to investigate the distribution of averages of 40 exponentials. Some parameters have been specified for us:

- Simulate the exponential distribution with “`rexp(n, lambda)`”
- Use  $\lambda = 0.2$  for all simulations
- The mean of the exponential distribution is  $1/\lambda$
- The standard deviation of the exponential distribution is  $1/\lambda$
- Do 1000 simulations ( $n=1000$ ) of 40 exponentials

First, we generate a population of exponentials. We do this using the R code provided for exponentials, and the values provided for  $\lambda$  and  $n$ .

```
pop <- rexp(1000, .2)
```

Next we generate samples of 40 exponentials, 1000 times, and take the mean of each. This is our simulation of multiple samples of the larger population.

```
sample.means <- NULL
for(i in 1:1000) sample.means = c(sample.means, mean(rexp(40, .2)))
```

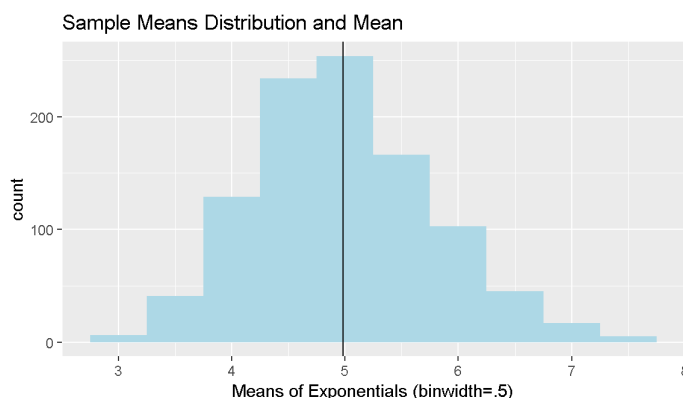
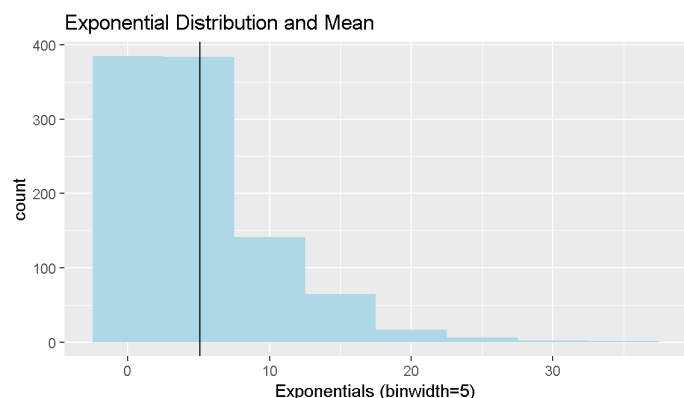
### Sample Mean versus Theoretical Mean

Per the CLT, the distribution of averages/means of the samples will approach normal, with a center/mean of the population mean.

We were given the theoretical mean of  $1/\lambda$ , which works out to  $1/.2$  or 5. Let's see how this compares to the mean of the populations we simulated, and the mean of the sample means.

```
## Population (1000) mean = 5.08529936893057
## Means of 1000 Samples of 40 = 4.98521953431527
## Theoretical mean = 5
```

Let's visualize the distributions with their means.



We can see that both the means are very close to the theoretical mean of 5.

### Sample Variance versus Theoretical Variance

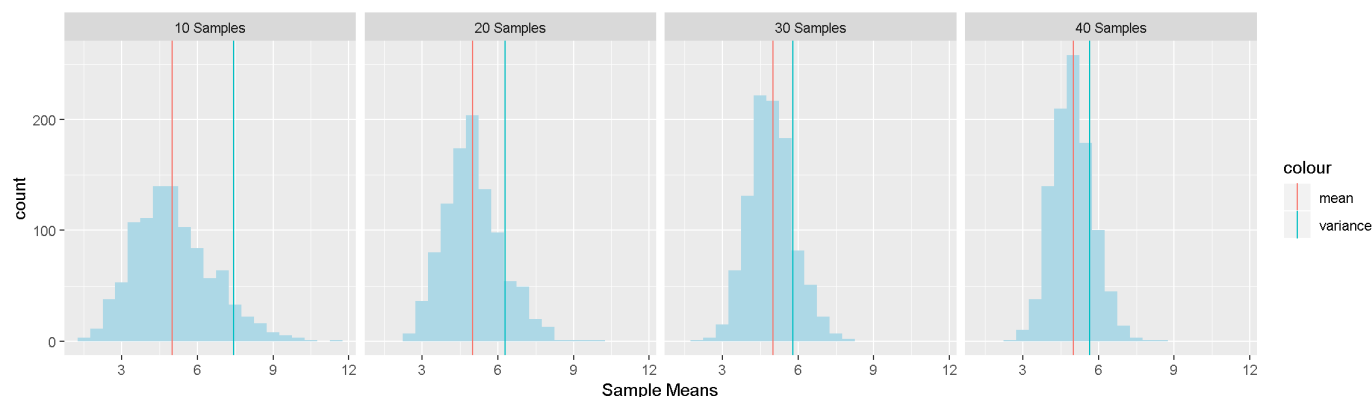
We know that the standard deviation =  $\sqrt{\text{variance}}$  and the theoretical  $\text{sd} = 1/\lambda = 5$ . Thus, the theoretical variance is 25.

The CLT states that the variance of the sample means will approach the theoretical variance /  $n$ . With our sample size of 40, this gives

25/40 = .625.

If we consider 4 different sets of sample means which vary by the number of samples in the mean (n), we see the variance (which is the sd squared), decline from the theoretical variance of 25.

```
## 10 Samples Variance: 2.43700615371867
## 20 Samples Variance: 1.29571436979027
## 30 Samples Variance: 0.799612814435421
## 40 Samples Variance: 0.641248583883752
```



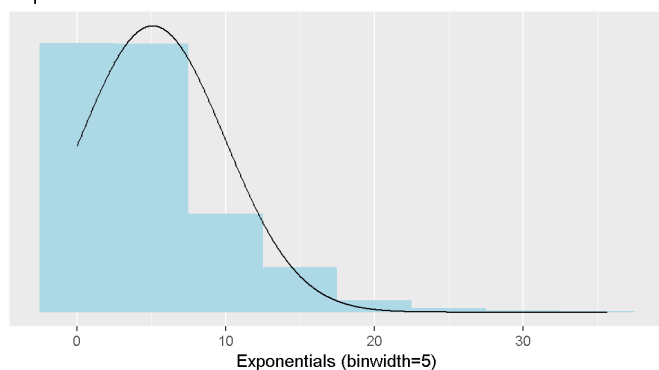
We have printed out the variance per sample size and visualized in a plot. With increasing sample size, we see the spread (variance) of the normal distribution decrease as the distribution clusters more tightly around the mean.

The actual variance of our 40 samples is 0.641 which is close to the theoretical variance of 25/40 or .625.

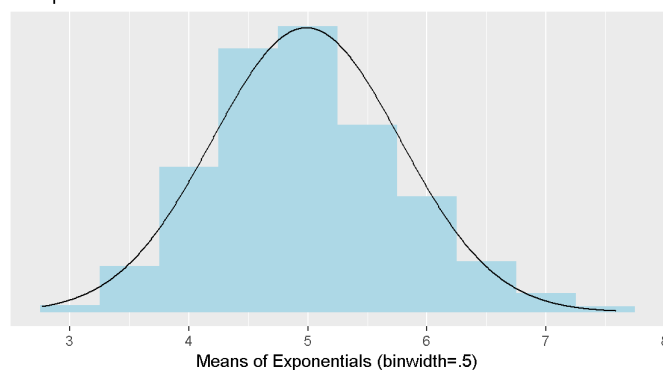
## Distribution

We will compare the simulated populations of exponentials and the sample means to a normal curve by plotting a normal on top of each. There are other methods that could be used, such as the Shapiro-Wilk test, which are out of the scope of this class.

Exponential Distribution



Sample Means Distribution



For the exponentials, there is no left tail and we see a distinct left skew. For the sample means, we see a clearly defined hump in the middle with tails on either side, indicative of a normal curve.

## Part 2: Basic Inferential Data Analysis

For part 2, we will analyze the ToothGrowth data in the R datasets package. This data set shows the effect of Vitamin C on tooth growth in guinea pigs.

## Exploratory Data Analysis

First we load the ToothGrowth data and perform some basic exploratory data analyses.

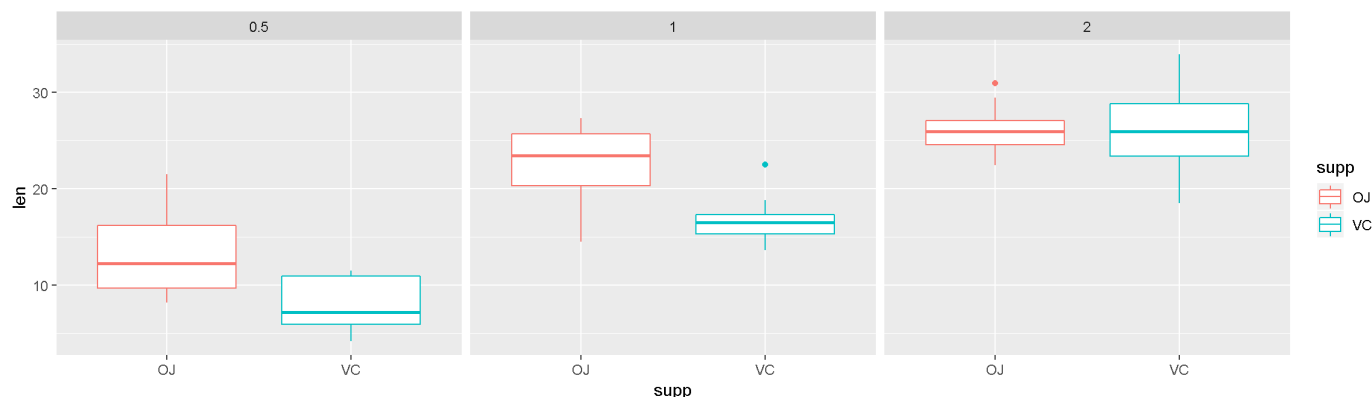
```
## 'data.frame': 60 obs. of 3 variables:
## $ len : num 4.2 11.5 7.3 5.8 6.4 10 11.2 11.2 5.2 7 ...
## $ supp: Factor w/ 2 levels "OJ","VC": 2 2 2 2 2 2 2 2 2 2 ...
## $ dose: num 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 0.5 ...
```

```
##
##      OJ VC
## 0.5 10 10
## 1   10 10
## 2   10 10
```

We see that there are 60 observations of 3 variables. There are 2 levels of supp and 3 distinct values of dose. There are 10 observations for each combination of supp and dose.

## Compare Tooth Growth

We will use confidence intervals and hypothesis testing (p-values) to compare tooth growth by supp and dose. Let's start with a visual of the data. We will do a hypothesis test to determine if there is significant difference between the two supp types at each dose level.



The null hypothesis ( $H_0$ ) will be that  $\text{mean}(\text{OJ}) = \text{mean}(\text{VC})$  for each dose. The alternative hypothesis ( $H_a$ ) will be that  $\text{mean}(\text{OJ}) > \text{mean}(\text{VC})$  for each dose.

Using the confidence interval calculation for comparing groups with unequal variances, means we will need to calculate the pooled standard error and the degrees of freedom. We'll compare the manual calculation with the R `t.test` command.

```
## [1] "Manual Calculation Confidence Interval: 2.34604034665748 , 8.15395965334252"
```

```
##
## Welch Two Sample t-test
##
## data:  ToothGrowth %>% filter(supp == "OJ", dose == "0.5") %>% pull(len) and ToothGrowth %>% filter(supp
== "VC", dose == "0.5") %>% pull(len)
## t = 3.1697, df = 14.969, p-value = 0.003179
## alternative hypothesis: true difference in means is greater than 0
## 95 percent confidence interval:
##  2.34604      Inf
## sample estimates:
## mean of x mean of y
##    13.23    7.98
```

Let's quickly look at the confidence intervals for the other doses.

```
## [1] "Dose = 1 Confidence Interval: c(3.35615763634793, Inf) / p-value: 0.000519187936149939"
```

```
## [1] "Dose = 2 Confidence Interval: c(-3.13349957439956, Inf) / p-value: 0.518074205638314"
```

## Conclusion

At dose = .5, there is a significant difference between OJ and VC. We know this because the 95% confidence interval does not include 0 (meaning it's less than 5% likely to have the means be equal), and the p-value is .003, which is less than .05.

At dose = 1, there is a significant difference between OJ and VC. Again, the confidence interval does not include 0, and the p-value is .0005 (well under .05).

At dose = 2, we can not reject the null hypothesis. The confidence interval includes 0, and the p-value is > .05.

These calculations are evidence of what we visually see in the plot of the differences between OJ and VC.

## Appendix

### Code

```
library(tidyverse)
library(ggplot2)
library(gridExtra)
```

```

# Part 1: Simulation Exercise

## Overview

## Simulations

pop <- rexp(1000, .2)
lg.pop <- rexp(2000, .2)

# combine into a data frame for plotting in ggplot2
df.pop <- rbind(as.data.frame(cbind(Exponential=pop, Group="1,000 Sample")),
               as.data.frame(cbind(Exponential=lg.pop, Group="2,000 Sample")))

ggplot(df.pop, aes(as.numeric(as.character(Exponential)))) +
  geom_histogram(binwidth = 5) +
  facet_grid(.~Group) +
  xlab("Exponentials (binwidth=5)")

sample.means <- NULL
for(i in 1:1000) sample.means = c(sample.means, mean(rexp(40, .2)))

### Sample Mean versus Theoretical Mean

cat(paste(" Population (1000) mean = ", mean(pop), "\n",
          "Means of 1000 Samples of 40 = ", mean(sample.means), "\n",
          "Theoretical mean = ", 1/.2))

p1 <- ggplot(as.data.frame(pop), aes(pop)) +
  geom_histogram(binwidth = 5, fill="lightblue") +
  xlab("Exponentials (binwidth=5)") +
  ggtitle("Exponential Distribution and Mean") +
  geom_vline(xintercept = mean(pop))

p2 <- ggplot(as.data.frame(sample.means), aes(sample.means)) +
  geom_histogram(binwidth = .5, fill="lightblue") +
  xlab("Means of Exponentials (binwidth=.5)") +
  ggtitle("Sample Means Distribution and Mean") +
  geom_vline(xintercept = mean(sample.means))

grid.arrange(p1, p2, ncol=2)

### Sample Variance versus Theoretical Variance

sample.means2 <- NULL
c(for(i in 1:1000) sample.means2 = c(sample.means2, mean(rexp(10, .2))),
  for(i in 1:1000) sample.means2 = c(sample.means2, mean(rexp(20, .2))),
  for(i in 1:1000) sample.means2 = c(sample.means2, mean(rexp(30, .2))),
  for(i in 1:1000) sample.means2 = c(sample.means2, mean(rexp(40, .2))))

# combine into a data frame for plotting in ggplot2
df.sample.means <- rbind(as.data.frame(cbind(Sample=sample.means2[1:1000], Group="10 Samples")),
                        as.data.frame(cbind(Sample=sample.means2[1001:2000], Group="20 Samples")),
                        as.data.frame(cbind(Sample=sample.means2[2001:3000], Group="30 Samples")),
                        as.data.frame(cbind(Sample=sample.means2[3001:4000], Group="40 Samples")))

df.intercepts <- as.data.frame(cbind(Variance= c(var(sample.means2[1:1000]),
                                                var(sample.means2[1001:2000]),
                                                var(sample.means2[2001:3000]),
                                                var(sample.means2[3001:4000])),
                                Group=c("10 Samples", "20 Samples", "30 Samples", "40 Samples")))

cat(paste(" 10 Samples Variance: ",
          df.intercepts[1,1], "\n",
          "20 Samples Variance: ",
          df.intercepts[2,1], "\n",
          "30 Samples Variance: ",
          df.intercepts[3,1], "\n",
          "40 Samples Variance: ",
          df.intercepts[4,1]))

ggplot(df.sample.means, aes(x=as.numeric(as.character(Sample)))) +
  geom_histogram(binwidth = .5, fill="lightblue") +

```

```

geom_vline(data = df.intercepts, aes(xintercept = 5+as.numeric(as.character(Variance))), colour="variance")
) +
geom_vline(aes(xintercept = 5, colour="mean")) +
facet_grid(.~Group) +
xlab("Sample Means")

### Distribution

p1 <- ggplot(as.data.frame(pop), aes(pop)) +
  geom_histogram(binwidth = 5, fill="lightblue", aes(y=..density..)) +
  xlab("Exponentials (binwidth=5)") +
  ggtitle("Exponential Distribution") +
  stat_function(fun = dnorm, n = 1000, args = list(mean = mean(pop), sd = sd(pop))) +
  ylab("") +
  scale_y_continuous(breaks = NULL)

p2 <- ggplot(as.data.frame(sample.means), aes(sample.means)) +
  geom_histogram(binwidth = .5, fill="lightblue", aes(y=..density..)) +
  xlab("Means of Exponentials (binwidth=.5)") +
  ggtitle("Sample Means Distribution") +
  stat_function(fun = dnorm, n = 1000, args = list(mean = mean(sample.means), sd = sd(sample.means))) +
  ylab("") +
  scale_y_continuous(breaks = NULL)

grid.arrange(p1, p2, ncol=2)

# Part 2: Basic Inferential Data Analysis

## Exploratory Data Analysis

# load data
data(ToothGrowth)
str(ToothGrowth)
table(ToothGrowth$dose, ToothGrowth$supp)

## Compare Tooth Growth

ggplot(ToothGrowth, aes(supp, len)) +
  geom_boxplot(aes(colour=supp)) +
  facet_grid(.~dose)

mean.OJ.0.5 <- mean(ToothGrowth%>%filter(supp=="OJ",dose=="0.5")%>%pull(len))
mean.VC.0.5 <- mean(ToothGrowth%>%filter(supp=="VC",dose=="0.5")%>%pull(len))
sd.OJ.0.5 <- sd(ToothGrowth%>%filter(supp=="OJ",dose=="0.5")%>%pull(len))
sd.VC.0.5 <- sd(ToothGrowth%>%filter(supp=="VC",dose=="0.5")%>%pull(len))

sp = sqrt((sd.OJ.0.5^2)/10 + (sd.VC.0.5^2)/10) # pooled standard error
sx = (sd.OJ.0.5^2)/10
sy = (sd.VC.0.5^2)/10
df = (sx+sy)^2/(sx^2/9 + sy^2/9) # degrees of freedom

paste("Manual Calculation Confidence Interval:",
      (mean.OJ.0.5-mean.VC.0.5)+-1*qt(.95,df)*sp, # confidence interval for unequal variances
      ", ",
      (mean.OJ.0.5-mean.VC.0.5)+qt(.95,df)*sp)

t.test(ToothGrowth%>%filter(supp=="OJ", dose=="0.5")%>%pull(len),
      ToothGrowth%>%filter(supp=="VC", dose=="0.5")%>%pull(len),
      paired=FALSE, var.equal=FALSE, alternative="greater", conf.level=.95) # t test for unequal variances

paste("Dose = 1 Confidence Interval:",
      t.test(ToothGrowth%>%filter(supp=="OJ", dose=="1")%>%pull(len),
            ToothGrowth%>%filter(supp=="VC", dose=="1")%>%pull(len),
            paired=FALSE, var.equal=FALSE, alternative="greater", conf.level=.95)[4],
      " / p-value:",
      t.test(ToothGrowth%>%filter(supp=="OJ", dose=="1")%>%pull(len),
            ToothGrowth%>%filter(supp=="VC", dose=="1")%>%pull(len),
            paired=FALSE, var.equal=FALSE, alternative="greater", conf.level=.95)[3]) # t test for dose = 1

paste("Dose = 2 Confidence Interval:",
      t.test(ToothGrowth%>%filter(supp=="OJ", dose=="2")%>%pull(len),
            ToothGrowth%>%filter(supp=="VC", dose=="2")%>%pull(len),
            paired=FALSE, var.equal=FALSE, alternative="greater", conf.level=.95)[4],

```

```
" / p-value:",
t.test(ToothGrowth%>%filter(supp=="OJ", dose=="2")%>%pull(len),
       ToothGrowth%>%filter(supp=="VC", dose=="2")%>%pull(len),
       paired=FALSE, var.equal=FALSE, alternative="greater", conf.level=.95)[3])# t test for dose = 2

## Conclusion
```