

# Statistical Inference - Simulation Exercise

*Abdul Rasheed Narejo*

*18/09/2018*

```
# load libraries
library(ggplot2) # for visualization
```

## A simulation exercise

In this project you will investigate the exponential distribution in R and compare it with the Central Limit Theorem. The exponential distribution can be simulated in R with `rexp(n, lambda)` where `lambda` is the rate parameter. The mean of exponential distribution is  $1/\lambda$  and the standard deviation is also  $1/\lambda$ . Set `lambda = 0.2` for all of the simulations. You will investigate the distribution of averages of 40 exponentials. Note that you will need to do a thousand simulations.

## Simulation

set the variables

```
# define the variables
n = 40
lambda = 0.2
simulations = 1000
```

set seed and generate the data

```
set.seed(1234)
data <- matrix(rexp(n = n * simulations, rate = lambda), nrow = simulations)
```

calculate mean of simulations

```
# calculate mean of each simulation
mean.data <- rowMeans(data)
```

## Comparing sample statistics with theoretical

Calculate mean, variance and standard deviation for sample generated from simulation

```
# sample mean
sam.mean <- round(mean(mean.data), 5)
sam.mean
```

```
## [1] 4.97424
```

```
# sample variance
sam.var <- round(var(mean.data), 5)
sam.var
```

```
## [1] 0.59497
```

```
# sample standard deviation
sam.sd <- sam.var^0.5
sam.sd
```

```
## [1] 0.771343
```

Calculate theoretical mean, variance and standard deviation

```
# theoretical mean
theo.mean <- 1 / lambda
theo.mean
```

```
## [1] 5
```

```
# theoretical variance
theo.var <- (1/lambda)^2/n
theo.var
```

```
## [1] 0.625
```

```
# theoretical standard deviation
theo.sd <- theo.var ^ 0.5
theo.sd
```

```
## [1] 0.7905694
```

Samople vs theoretical

- Sample mean is 4.97424 compared to theoretical mean of 5, indicating a negligible difference of -0.52%.
- Sample variance is 0.59497 compared to theoretical mean of 0.625, indicating a negligible difference of -4.8%.
- Sample standard deviation is 0.771343 compared to theoretical standard deviation of 0.7905694, indicating a negligible difference of -2.43%.

```
mean.data <- data.frame(mean.data)
names(mean.data) <- c("mean")

plot <- ggplot(mean.data, aes(x = mean))
plot <- plot + geom_histogram(aes(y=..density..), breaks=seq(2,8, by = .1),
                             col = "dark blue", fill = "dark blue", alpha = .25)
plot <- plot + labs(title = "Mean of exponential distribution",
                   x = "Average mean of 40 exponentials", y="Density")

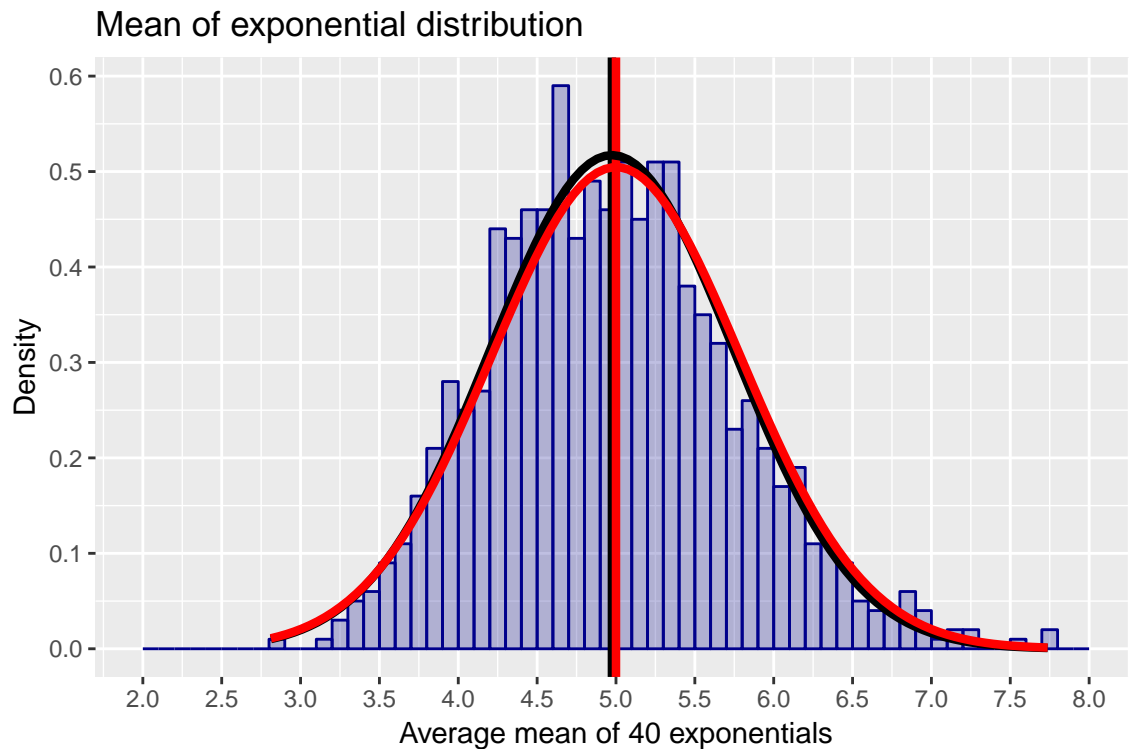
plot <- plot + scale_x_continuous(breaks = seq(2,8, by=0.5))
plot <- plot + scale_y_continuous(breaks = seq(0,0.6, by=0.1))

plot <- plot + geom_vline(aes(xintercept = sam.mean), col = "black",
                          size=1.5)
plot <- plot + geom_vline(aes(xintercept = 5), col = "red", size=1.5)
```

```
plot <- plot + stat_function(fun = dnorm, args = list(mean = sam.mean, sd = sam.sd),
  color = "black", size = 1.5)

plot <- plot + stat_function(fun = dnorm, args = list(mean = theo.mean, sd = theo.sd),
  color = "red", size = 1.5)

plot
```



Simulation mean (black curve) and theoretical mean (red curve) nearly overlap perfectly and depict approximately normal distributions over the dataset.

### Simulation vs. theoretical confidence interval

```
samCI <- round(sam.mean + c(-1,1)*1.96*sam.sd/sqrt(n),3)
theoCI <- round(theo.mean + c(-1,1)*1.96*theo.sd/sqrt(n),3)
```

95% confidence interval for simulation is (4.735, 5.213) while 95% confidence interval for theoretical distribution is (4.755, 5.245).

### Conclusion

As shown above, the mean and variance of 40 randomly generated sample values of exponential distribution are very close to the theoretical mean and variance. Hence, the sample distribution is approximately normal as proven by mean, sd, variance and confidence interval.

## Appendix

### What is exponential distribution ?

In probability theory and statistics, the exponential distribution (also known as the negative exponential distribution) is the probability distribution that describes the time between events in a Poisson point process, i.e., a process in which events occur continuously and independently at a constant average rate.

The exponential distribution may be useful to model events such as;

- The time between goals scored in a World Cup soccer match
- The duration of a phone call to a help center
- The time between meteors greater than 1 meter diameter striking earth
- The time between successive failures of a machine
- The time from diagnosis until death in patients with metastatic cancer
- The distance between successive breaks in a pipeline

### QQ Plot

```
ggplot(mean.data, aes(sample = mean)) +  
  stat_qq() +  
  labs(title = "Quartile plot of Sample data mean",  
        y="sample") +  
  
  stat_qq_line(colour="red", size=1.5)
```

