

Exponential Distribution and Central Limit Theorem

Ivana Lipnerová

December 20, 2015

Study

Overview

The goal of this project is to investigate exponential distribution compared to distribution of averages of 40 exponentials, as these according to Central Limit Theorem should approximate normal distribution. I will test this via 1000 simulations of drawing 40 random exponentials with set parametr lambda ($\lambda = 0.2$). In the instruction there is useful information about the mean of exponential distribution (mean = $1/\lambda$) and the standard deviation of same distribution (sd = $1/\lambda$).

Simulations

```
#load necessary libraries:
library(ggplot2); library(gridExtra); library(knitr)

# Set lambda, number of simulations and drawn samples:
lambda <- 0.2
n_simulations <- 1000
n <- 40

# Set random seed, so that the conclusion will be valid for every run of knitr:
set.seed(12345)

# Create random exponential variables:
sim_vars <- matrix(data=rexp(n*n_simulations, rate=lambda), nrow=n_simulations, ncol=n)

# Create means:
sim_means <- rowMeans(sim_vars)
```

Sample Mean versus Theoretical Mean

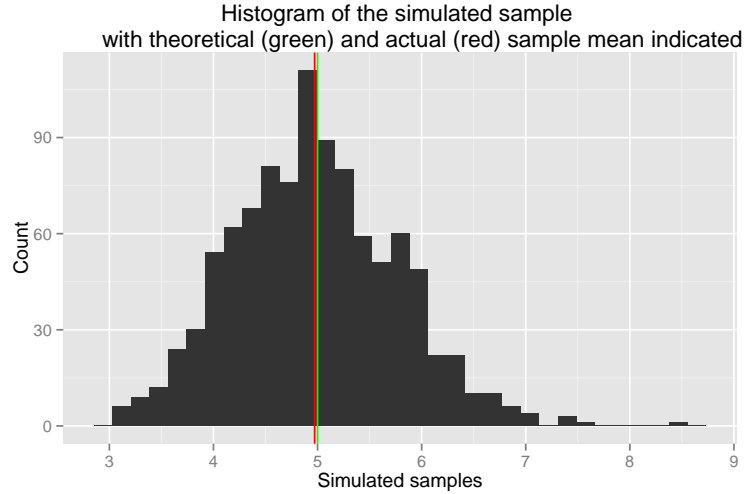
The mean of simulated averages can be safely computed by base function mean(). As was stated in overview, the theoretical mean of exponential distribution is $1/\lambda$, with $\lambda = 0.2$.

```
means<-data.frame("Mean"=c(mean(sim_means), 1/lambda),
                  row.names=c("Sample mean", "Theoretical mean"))
kable(x=means, digits=3, align="c", caption="Table of sample and
corresponding theoretical mean")
```

Table 1: Table of sample and corresponding theoretical mean

	Mean
Sample mean	4.972
Theoretical mean	5.000

As we can see, the actual sample mean is not so far from the theoretically predicted one. Let's see how it stands in histogram of simulated data (Fig.1., see Appendix for code).



Sample Variance versus Theoretical Variance

The variance of simulated averages can be safely computed by base function `var()`. The theoretical standard deviation of exponential distribution is $1/\lambda$. Given the Central Limit Theorem, the theoretical variance is

$$\frac{\sigma^2}{n} = \frac{(1/\lambda)^2}{n}$$

```
var_sim <- var(sim_means)
var_theo <- ((1/lambda)^2)/n

vars<-data.frame("Variances"=c(var_sim, var_theo),
                 row.names=c("Sample variance", "Theoretical variance"))
kable(x=vars, digits=3, align="c", caption="Table of sample
and corresponding theoretical variance")
```

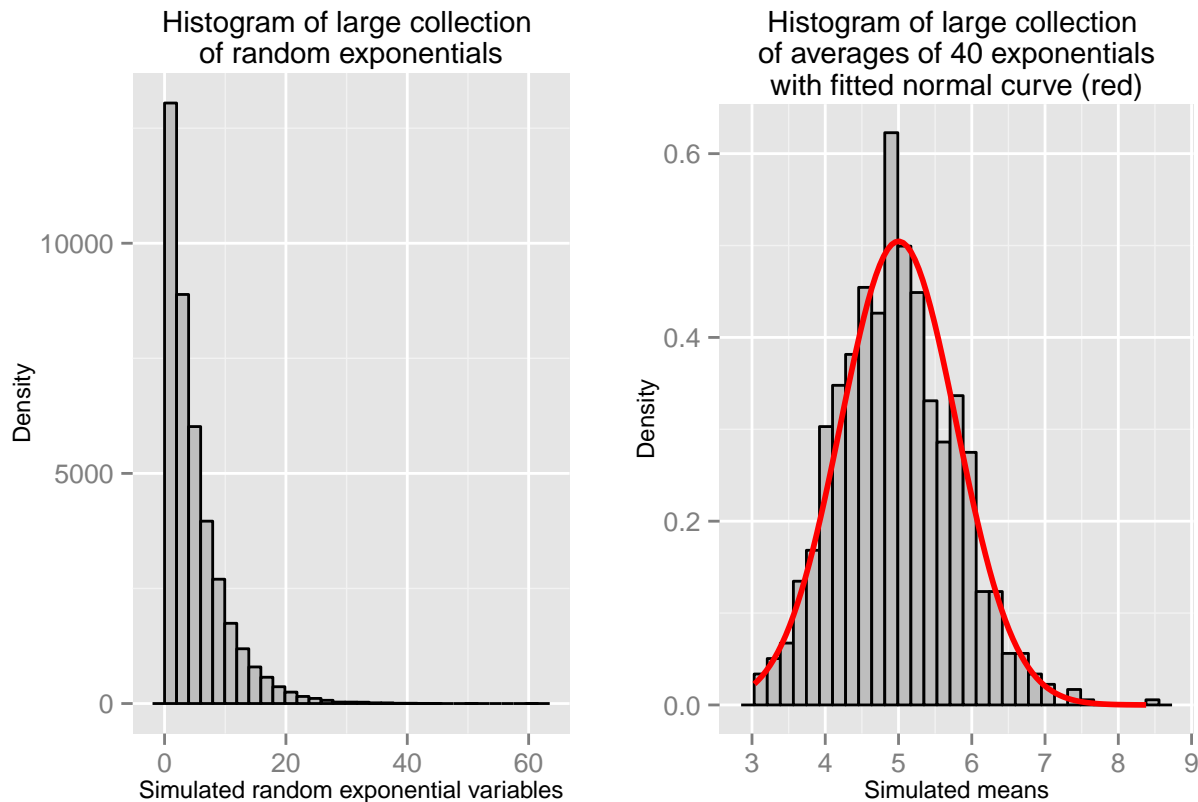
Table 2: Table of sample and corresponding theoretical variance

	Variances
Sample variance	0.616
Theoretical variance	0.625

Distribution

I decided to look into this question by showing the fit of normal curve to the histogram (see plot below) and by testing the normality of averages by Shapiro-Wilk test of normality (at the bottom of the document).

For the first part, there are two histograms. The left one shows the original simulated data (thus showing them as exponentials), the right one shows the drawn averages of 40 exponentials with fitted normal curve for theoretical normal distribution given set lambda (Fig. 2., see Appendix for code).



The plots are illustrative, though at least simple test for normality should be done. In Shapiro-Wilk test of normality the p-value lower than 0.01 indicates non-normal distribution of tested data.

```
shapiro.test(sim_means)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  sim_means  
## W = 0.99438, p-value = 0.0008625
```

Conclusion

Thus I would conclude that the tested collection of averages of 40 exponentials is just approximately normal. Given the lambda and setted random seed, to achieve normality of averages one should use larger dataset (e.g. simulate more data).

Appendix

This appendix contains code used for generating plots.

Fig. 1. Histogram of the simulated sample with theoretical (green) and actual (red) sample mean indicated:

```
ggplot(as.data.frame(sim_means), aes(sim_means)) +  
  geom_histogram() +  
  geom_vline(xintercept=mean(sim_means), colour="red") +  
  geom_vline(xintercept=1/lambda, colour="green") +  
  labs(title="Histogram of the simulated sample  
    with theoretical (green) and actual (red) sample mean indicated",  
    x="Simulated samples", y="Count")
```

Fig. 2. Two histograms, left one of raw exponential data, right one of drawn averages.

```
# Plot of raw simulated data:  
p1 <- ggplot(as.data.frame(as.vector(sim_vars)), aes(as.vector(sim_vars))) +  
  geom_histogram(alpha=0.5, fill="grey", col="black") +  
  ggtitle("Histogram of large collection  
of random exponentials") +  
  xlab("Simulated random exponential variables") +  
  ylab("Density") +  
  theme(title = element_text(size = rel(0.75), hjust = 0.5))  
  
# Plot of averages:  
p2 <- ggplot(as.data.frame(sim_means), aes(sim_means)) +  
  geom_histogram(aes(y=..density..), alpha=0.5, fill="grey", col="black") +  
  #that aes() necessary for normal curve displayed correctly  
  stat_function(fun=dnorm, colour="red", lwd=1,  
    args=list(mean = 1/lambda, sd = sqrt(var_theo))) +  
  #necessary manually set mean and sd for a curve  
  labs(title="Histogram of large collection  
of averages of 40 exponentials  
with fitted normal curve (red)",  
    x="Simulated means",  
    y="Density") +  
  theme(title = element_text(size = rel(0.75), hjust = 0.5))  
  
# Arrange them:  
grid.arrange(p1, p2, ncol=2, nrow=1)
```