

Simulation of a Convergence to the Normal Distribution

Jim Callahan

September 23, 2015

Overview

It is a useful concept in probability and statistics that even if a probability distribution is not normal (does not follow the “bell curve”) that if you calculate means (averages) of small subsamples from the distribution the distribution of the means will be normal, as long as the underlying distribution has a finite variance (for example, not a power law distribution) and meets a few less well known mathematical properties.

In this project we demonstrate how averages of subsets of an asymmetric distribution (in this case the **exponential distribution**) converge to a symmetric normal distribution.

The convergence of averages of subsets to normal distribution is an application of the **Central Limit Theorem (CLT)**. The “Statistical Inference” course presentation states: “For our purposes, the **CLT** states that the distribution of averages of **iid [independent, identically distributed]** variables, properly normalized, **becomes that of a standard normal as the sample size increases**” slide from **Dr. Brian Caffo’s** “*Statistical Inference*” class slide 8/20 “A Trip to Asymptopia” in Asymptopia.pdf

Simulations

This project uses the programming language **R** to simulate draws from an **exponential distribution**.

We first look graph and simulation of the **exponential distribution** and then compute and graph a mean (average) of subsets 40 observations each from the exponential distribution and finally compare the distribution of the means with the normal distribution to see if the process of applying the means has “shape shifted” the distribution from exponential to normal.

The draws from the **exponential distribution** are simulated in **R** with a random number function. Random number functions typically return values either a normal or a uniform distribution. This simulation depends on a family of random number functions in **R** that are modified to return random values from specific probability distributions other than the normal distribution.

Specifically, **rexp()** is in **R** a specialized random number function that returns values from the **exponential distribution**:

```
x <- rexp(n, lambda)
where:
n      = the number of draws
lamda  = the rate parameter
        (for purposes of this project lambda = .2)
x      = the values of x returned by the draw.
```

In **theoretical** terms, the **exponential distribution** has the property that the expected value of the **mean** and **standard deviation** (the standard deviation is the square root of the variance) are both equal to **the inverse of lambda (1/lambda)**.

Exponential Distribution

```
mean          = 1/lambda (in this project 1/.2 = 5)
std deviation  = 1/lambda (in this project 1/.2 = 5)
std deviation  = squareroot(variance)
variance       = square of the standard deviation (sd^2)
```

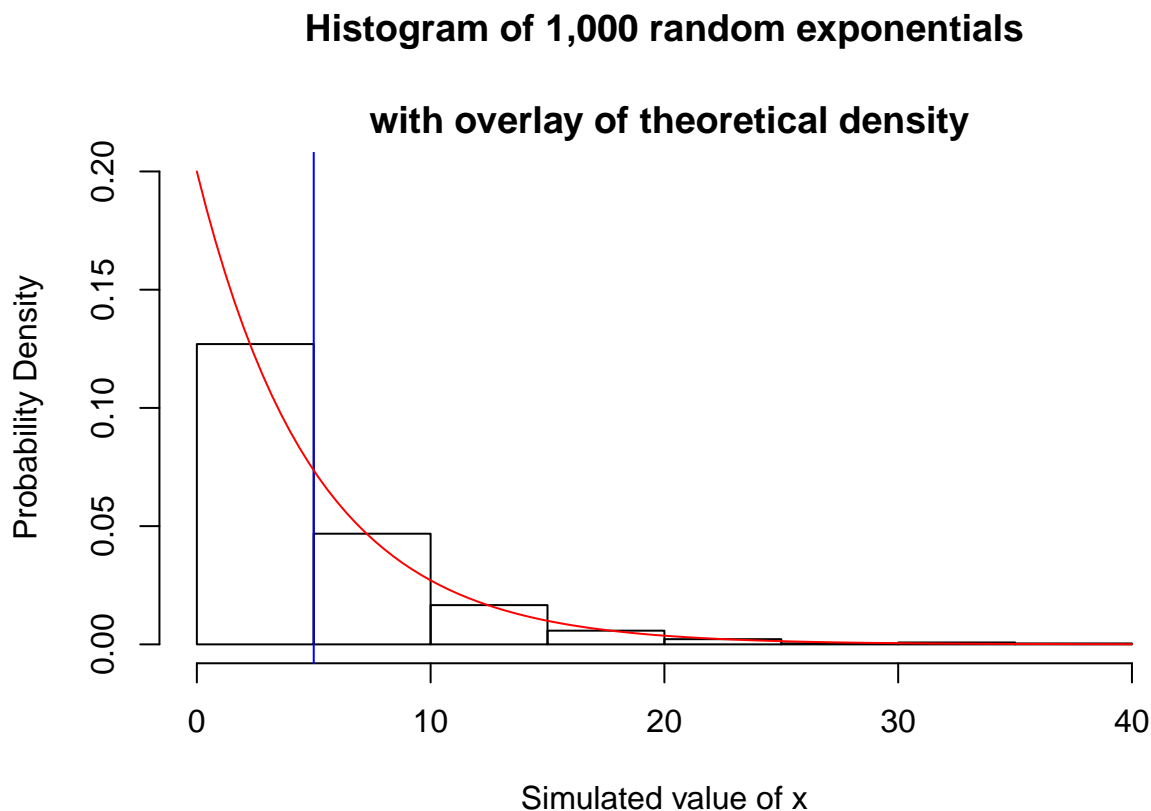
We begin with graph and simulation of the **exponential distribution** itself; after that we will compute and graph a mean (average) of sets 40 observations each from 1,000 simulations.

```
# set the random number seed so the results are reproducible.
set.seed(1234)
# Simulate 1000 random exponentials
x <- rexp(1000, rate = 0.2)

# Make the Y-axis start at zero and go to the maximum value
# in the range() function.
ylim <- range(0, .20)

# Plot the empirical results of the simulation
# Use density instead of frequency,
# so the theoretical curve will fit.
hist(x, freq = F, ylim = ylim,
     main = "Histogram of 1,000 random exponentials
           \n with overlay of theoretical density",
     xlab = "Simulated value of x",
     ylab = "Probability Density"
    )

# Overlay a plot a theoretical exponential density curve
abline(v = 5, col="Blue")
curve(dexp(x, rate = 0.2), add=T, col="Red")
```



A simulation of 1,000 values of the exponential distribution with a histogram of the raw (untransformed)

simulated values of x without taking an average. For comparison we overlaid (the smooth red curve) the theoretical exponential distribution. The vertical blue line shows the expected mean of 5 when $\lambda = .2$.

Sample Mean versus Theoretical Mean

We know the theoretical expected mean should approach 5 (when $\lambda = .2$); what do we get when we calculate the mean of the specific values of x we drew with the simulation?

```
mean(x)
```

```
## [1] 5.003067
```

Sample Variance versus Theoretical Variance

We also know the theoretical expected standard deviation should approach 5 (when $\lambda = .2$). Because the standard deviation is defined as the square root of the variance; the variance should be the square of the standard deviation; which in this case with a standard deviation of 5, we expect the variance (square of the standard deviation) to approach 25. When we calculate the standard deviation and variance of the specific values of x we drew with the simulation, we get:

```
# Theoretical variance should be the square of the standard deviation  
# The theoretical standard deviation equals 5; so the theoretical variance should equal 25  
sd(x)           # sample standard deviation
```

```
## [1] 5.056718
```

```
(sd(x))^2      # standard deviation squared
```

```
## [1] 25.5704
```

```
var(x)         # sample (empirical) variance
```

```
## [1] 25.5704
```

Distribution How one can tell the distribution is approximately normal?

```
# compare the distribution of # 1,000 random numbers  
# from the exponential probability distribution  
  
# Histogram template (from R help file for hist() function):  
# hist(x, breaks = "Sturges",  
#     freq = NULL, probability = !freq,  
#     include.lowest = TRUE, right = TRUE,  
#     density = NULL, angle = 45, col = NULL, border = NULL,  
#     main = paste("Histogram of" , xname),  
#     xlim = range(breaks), ylim = NULL,  
#     xlab = xname, ylab,  
#     axes = TRUE, plot = TRUE, labels = FALSE,  
#     nclass = NULL, warn.unused = TRUE, ...)
```

```

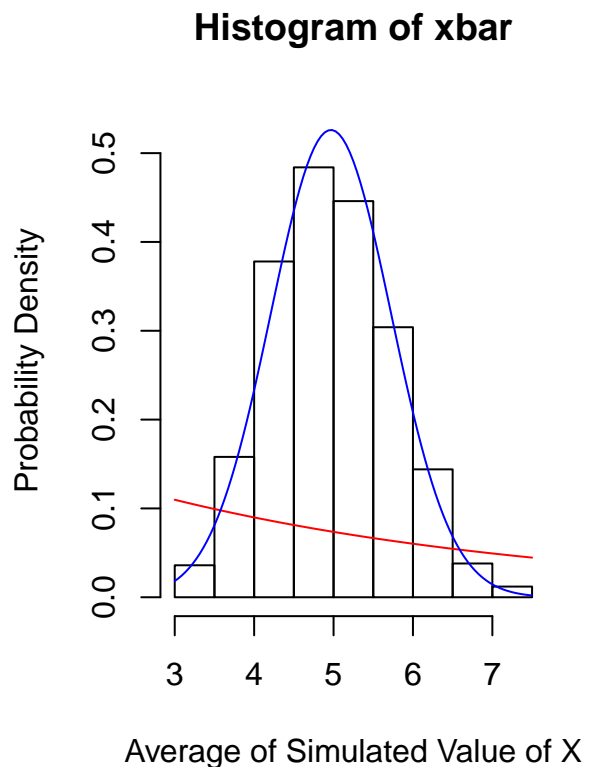
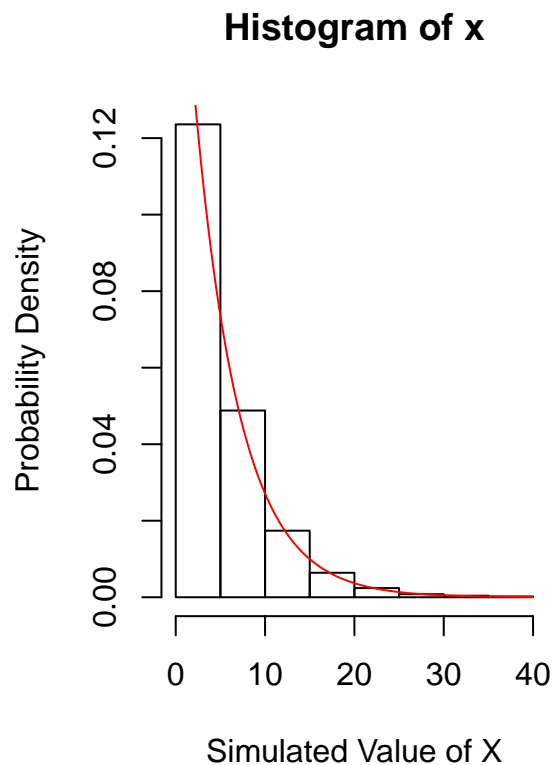
par(mfrow = c(1,2))

# Compare the distribution of 1000 random exponentials
# ylim <- .55 # (from h2 graph below)
x <- rexp(1000, rate = 0.2)
h1 <- hist(x, freq=F,
           xlab = "Simulated Value of X",
           ylab = "Probability Density" )
curve(dexp(x, rate = 0.2), add=T, col="Red") # Exponential Curve (Theoretical)

# with the distribution of 1000 averages of 40 random exponentials
xbar = NULL
for (i in 1 : 1000) xbar = c(xbar, mean(rexp(40, rate = 0.2)))
xsave <- x
x <- xbar
m <- mean(xbar)
s <- sd(xbar)
h2 <- hist(xbar, plot = F)

ylim = range(0, h2$density, max(h2$density)*1.1 )
hist(xbar, freq=F, ylim=ylim,
     xlab = "Average of Simulated Value of X",
     ylab = "Probability Density" )
curve(dexp(x, rate = 0.2), add=T, col="Red") # Exponential Curve (Theoretical)
curve(dnorm(x, m, s), add=T, col="Blue")    # Normal Curve (Theoretical)

```



Q-Q Plot with Sample Quantiles versus Theoretical Quantiles Recall that “the distribution of averages of iid [independent, identically distributed] variables, **properly normalized**, becomes that of a standard normal as the sample size increases.”

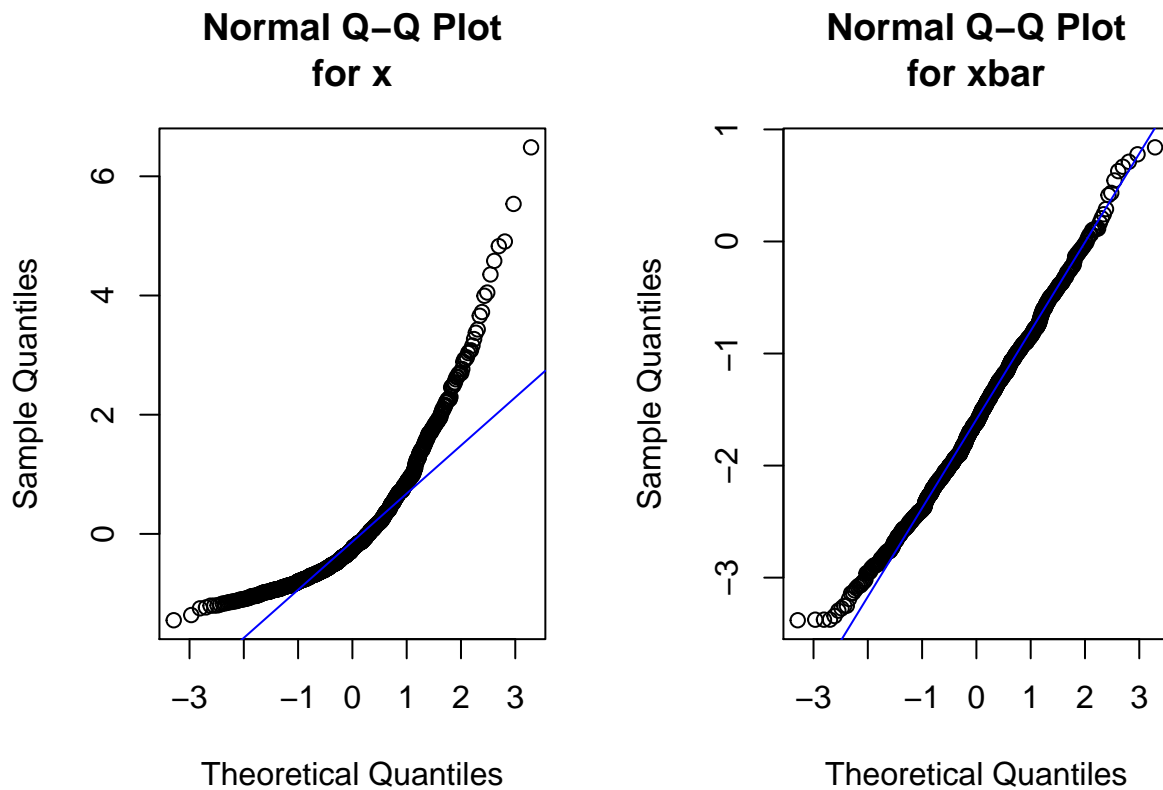
```
par(mfrow = c(1,2) )
x <- xsave

# Normalize x and x bar before plotting
normx    <- (x - xbar)/sd(x)
normxbar <- (xbar - mean(xbar)/sd(xbar))

# Q-Q Plot template (from R help file for qqnorm() function):
# qqnorm(y, ylim, main = "Normal Q-Q Plot",
#       xlab = "Theoretical Quantiles", ylab = "Sample Quantiles",
#       plot.it = TRUE, datax = FALSE, ...)

qqnorm(normx,    main = "Normal Q-Q Plot\nfor x")
qqline(normx,    col = "Blue")

qqnorm(normxbar, main = "Normal Q-Q Plot\nfor xbar")
qqline(normxbar, col = "Blue")
```



In both of the pair graphs the right hand graph shows that the mean of the 40 observations is closer to normal than the underlying exponential distribution.

References Statistical Inference Class

Dr. Brian Caffo, “*Statistical Inference*” class slides * Central Limit Theorem, slide 8/20 “A Trip to Asymptopia” in Asymptopia.pdf

External

Dr. Peter Dalggaard, “*Introductory Statistics with R*”, book * pages 31-32, histogram

* pages 64-65, Q-Q plots

Phillip K. Janert, “*Data Analysis with Open Source Tools*”, book * pages 25-26 “Optional: Comparing Distributions with Probability and QQ Plots”