

Statistical Inference - Coursera Week 3 Project: Sumilation Exercises

Girish Babu

21 September 2014

Goals of this analysis

The exponential distribution can be simulated in R with `rexp(n, lambda)` where `lambda` λ is the rate parameter. The mean of exponential distribution is $1/\lambda$ and the standard deviation is also $1/\lambda$. For this simulation, we set $\lambda = 0.2$. In this simulation, we investigate the distribution of averages of 40 numbers sampled from exponential distribution with $\lambda = 0.2$.

Perform a thousand simulated averages of 40 exponentials

```
set.seed(3)
lambda <- 0.2
sim_num <- 1000
sample_size <- 40
simulation <- matrix(rexp(sim_num * sample_size, rate=lambda), sim_num, sample_size)
means_of_row <- rowMeans(simulation)
```

Let us plot the the distribution of sample means.

```
# Histogram plot of averages
hist(means_of_row, breaks=50, prob=TRUE,
     main="Distribution of averages of samples,
     drawn from exponential distribution with lambda=0.2",
     xlab="")

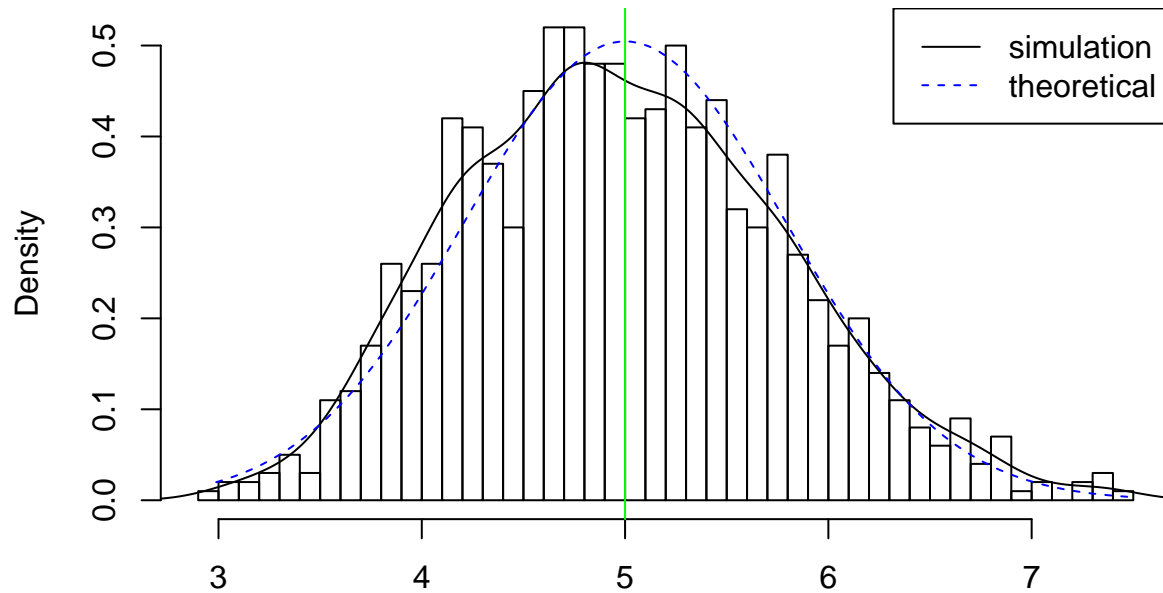
# Density of the averages of samples
lines(density(means_of_row))

# Theoretical center of distribution
abline(v=1/lambda, col="green")

# Theoretical density of the averages of samples
xfit <- seq(min(means_of_row), max(means_of_row), length=100)
yfit <- dnorm(xfit, mean=1/lambda, sd=(1/lambda/sqrt(sample_size)))
lines(xfit, yfit, pch=22, col="blue", lty=2)

legend('topright', c("simulation", "theoretical"), lty=c(1,2), col=c("black", "blue"))
```

Distribution of averages of samples, drawn from exponential distribution with lambda=0.2

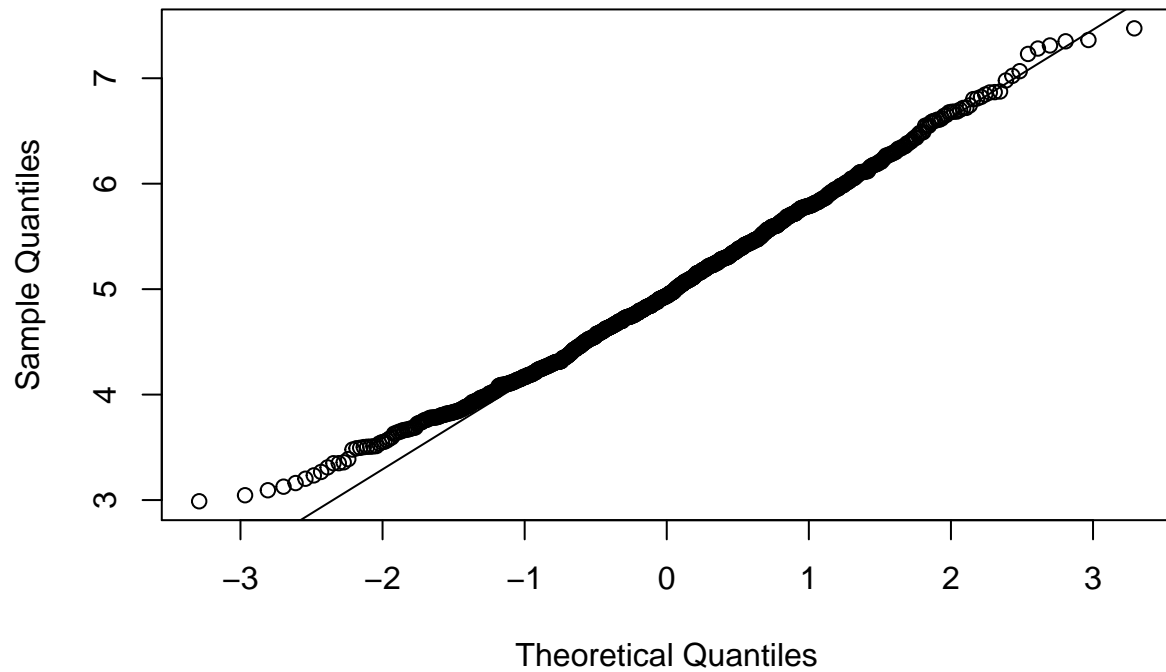


The distribution of sample means is centered at 4.9866 and the theoretical center of the distribution is $\lambda^{-1} = 5$. The variance of sample means is 0.6258 where the theoretical variance of the distribution is $\sigma^2/n = 1/(\lambda^2 n) = 1/(0.04 \times 40) = 0.625$.

Due to the Central Limit Theorem (CLT), the averages of samples follow normal distribution. The figure above also shows the density computed using the histogram and the normal density plotted with theoretical mean and variance values. Also, the q-q plot below suggests the normality.

```
qqnorm(means_of_row); qqline(means_of_row)
```

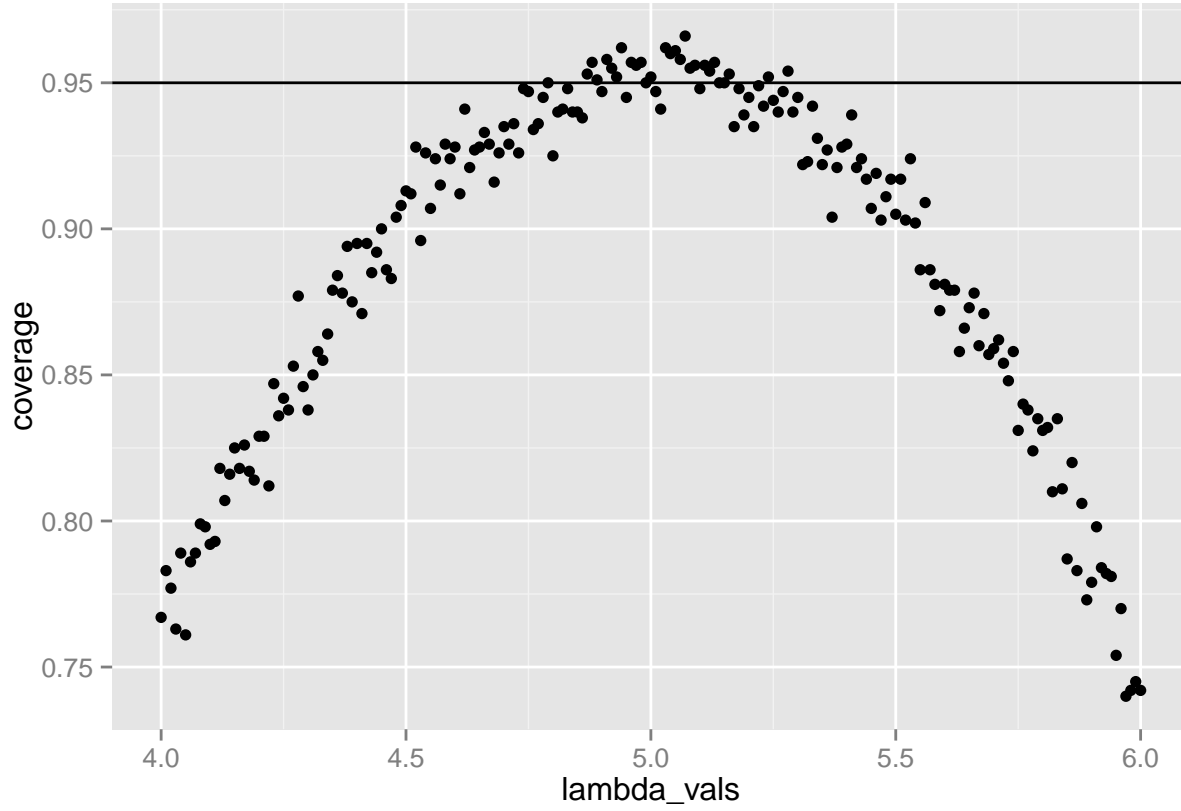
Normal Q-Q Plot



Finally, let's evaluate the coverage of the confidence interval for $1/\lambda = \bar{X} \pm 1.96 \frac{S}{\sqrt{n}}$

```
lambda_vals <- seq(4, 6, by=0.01)
coverage <- sapply(lambda_vals, function(lamb) {
  mu_hats <- rowMeans(matrix(rexp(sample_size * sim_num, rate=0.2),
                             sim_num, sample_size))
  ll <- mu_hats - qnorm(0.975) * sqrt(1/lambda**2/sample_size)
  ul <- mu_hats + qnorm(0.975) * sqrt(1/lambda**2/sample_size)
  mean(ll < lamb & ul > lamb)
})

library(ggplot2)
qplot(lambda_vals, coverage) + geom_hline(yintercept=0.95)
```



The 95% confidence intervals for the rate parameter (λ) to be estimated ($\hat{\lambda}$) are $\hat{\lambda}_{low} = \hat{\lambda}(1 - \frac{1.96}{\sqrt{n}})$ and $\hat{\lambda}_{upp} = \hat{\lambda}(1 + \frac{1.96}{\sqrt{n}})$. As can be seen from the plot above, for selection of $\hat{\lambda}$ around 5, the average of the sample mean falls within the confidence interval at least 95% of the time. Note that the true rate, λ is 5.