

Statistical Inference Project-1: Simulation

coolbhatt

Sunday, August 23, 2015

Statistical Inference Course Project, Part 1: Simulation Exercises

The exponential distribution can be simulated in R with `rexp(n, lambda)` where 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 $(\lambda=0.2)$.

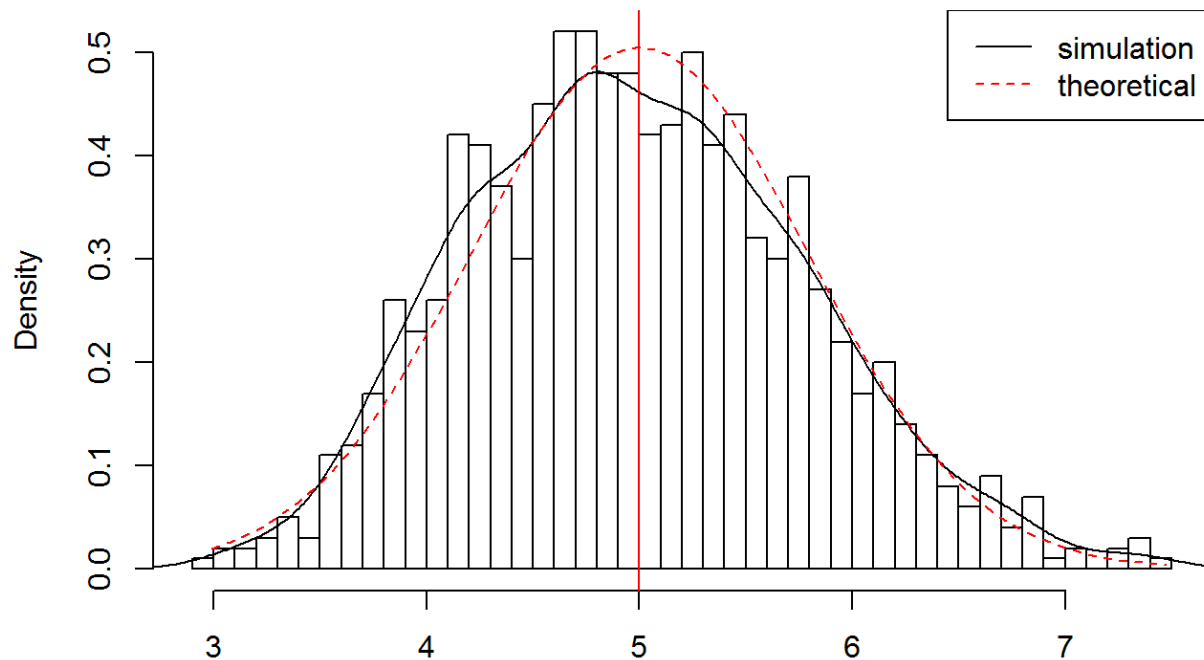
Let's do a thousand simulated averages of 40 exponentials.

```
set.seed(3)
lambda <- 0.2
num_sim <- 1000
sample_size <- 40
sim <- matrix(rexp(num_sim*sample_size, rate=lambda), num_sim, sample_size)
row_means <- rowMeans(sim)
```

The distribution of sample means is as follows.

```
# plot the histogram of averages
hist(row_means, 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(row_means))
# theoretical center of distribution
abline(v=1/lambda, col="red")
# theoretical density of the averages of samples
xfit <- seq(min(row_means), max(row_means), length=100)
yfit <- dnorm(xfit, mean=1/lambda, sd=(1/lambda/sqrt(sample_size)))
lines(xfit, yfit, pch=22, col="red", lty=2)
# add legend
legend('topright', c("simulation", "theoretical"), lty=c(1,2), col=c("black",
"red"))
```

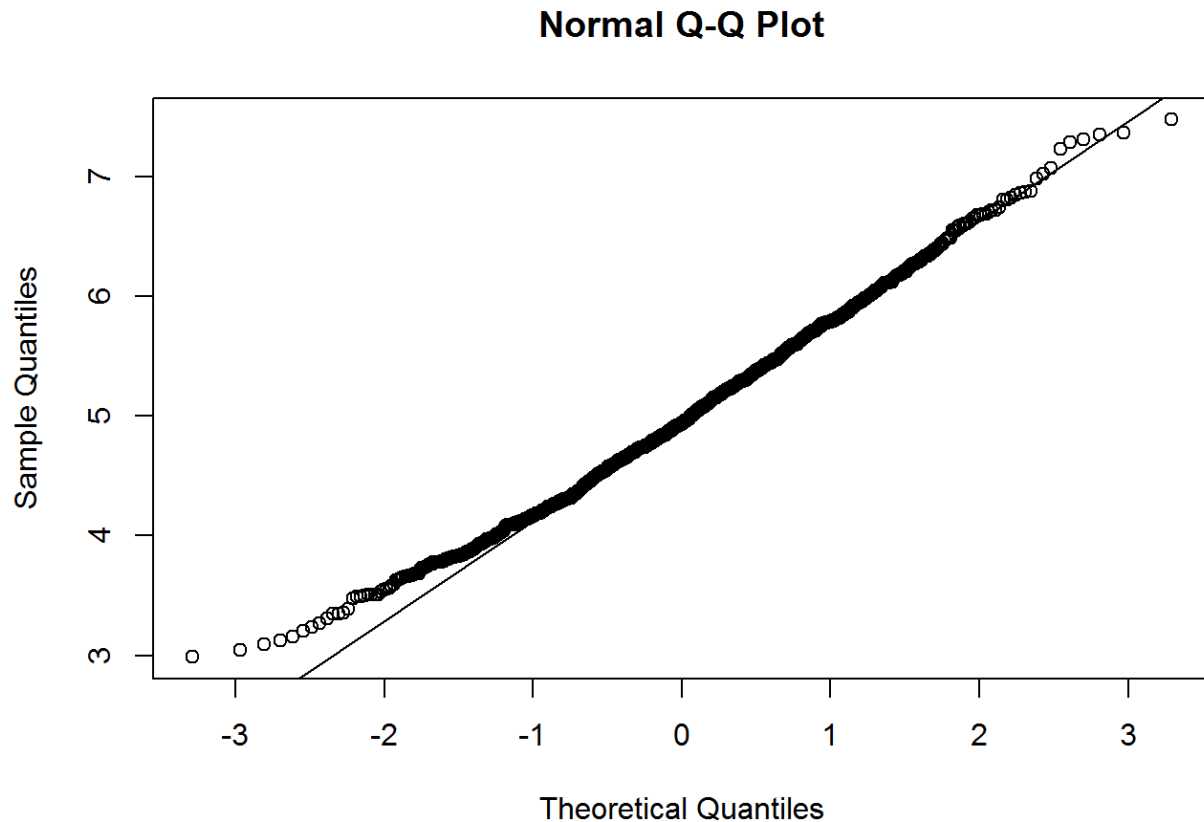
Distribution of averages of samples, drawn from exponential distribution with $\lambda=0.2$



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

Due to the central limit theorem, 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(row_means); qqline(row_means)
```



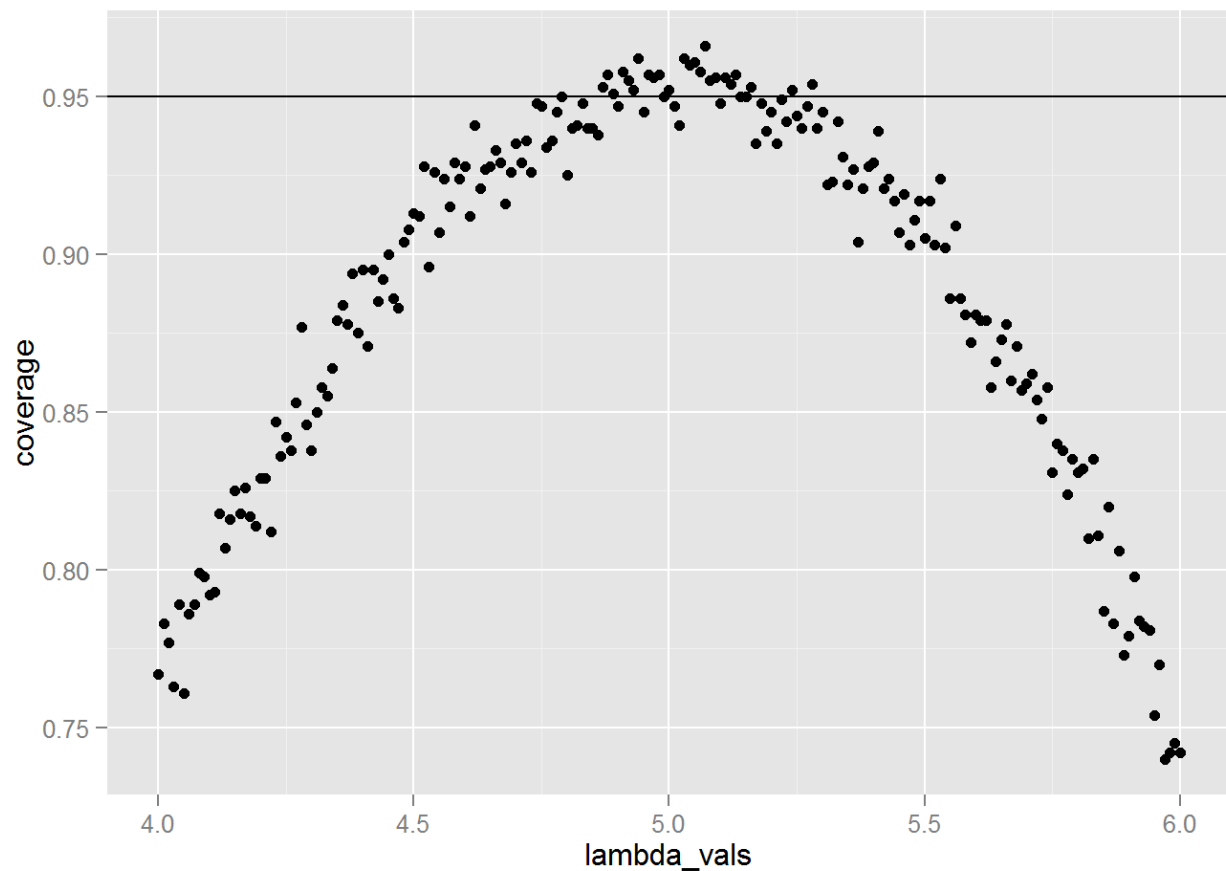
Finally, let's evaluate the coverage of the confidence interval for λ ($1/\lambda = \bar{X} \pm 1.96 \frac{S}{\sqrt{n}}$)

```
lambda_vals <- seq(4, 6, by=0.01)
coverage <- sapply(lambda_vals, function(lamb) {
  mu_hats <- rowMeans(matrix(rexp(sample_size*num_sim, rate=0.2),
                              num_sim, 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)
```

```
## Warning: package 'ggplot2' was built under R version 3.1.3
```

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



The 95% confidence intervals for the rate parameter (λ) to be estimated ($\hat{\lambda}$) are $\hat{\lambda}_{\text{low}} = \hat{\lambda}(1 - \frac{1.96}{\sqrt{n}})$ and $\hat{\lambda}_{\text{upp}} = \hat{\lambda}(1 + \frac{1.96}{\sqrt{n}})$. 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.