

Statistical Inference: Course Project

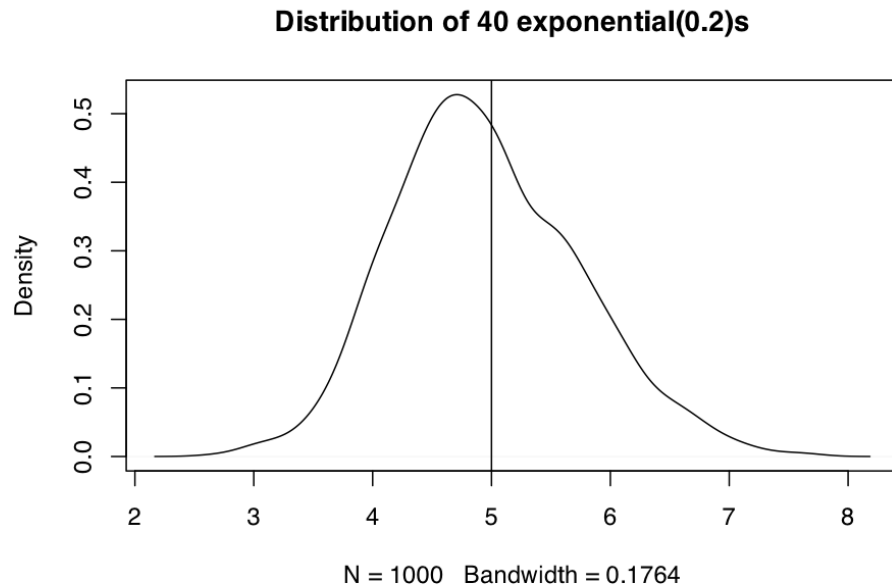
Peggy Fan

September 15, 2014

Part I. The purpose of this report is to compare the distribution of sample means of an exponential distribution with 1000 simulations and lambda at 0.2.

1. Where the distribution is centered at comparing to the theoretical center of the distribution.

I plotted the data of sample means (mean of 40 exponentials). The vertical ($x=5$) is the theoretical center of the distribution. The plot shows that the peak of the distribution of the sample means is a little less than 5.



2. How variable it is comparing to the theoretical variance of the distribution. The standard deviation (SD) of each sample of 40 is centered around the true SD, $1/\lambda$. But the SD of the means of the 1000 40-exponentials samples is the standard error, which is S/\sqrt{n} . So the theoretical standard deviation of the 1000 means should be $(1/\lambda)/\sqrt{40}$, and the variance should be

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

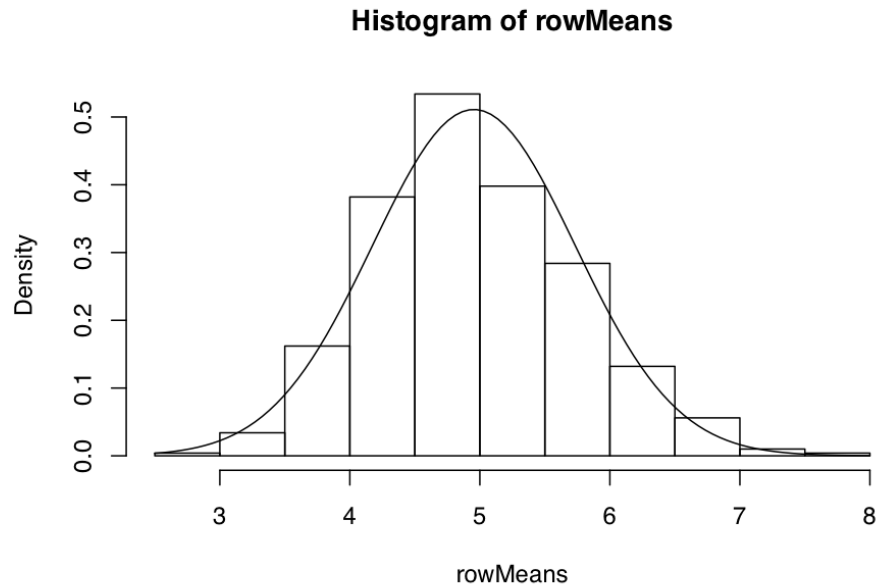
The variance of the dataset, by calculating its standard deviation then squaring it, gives

```
## [1] 0.6092
```

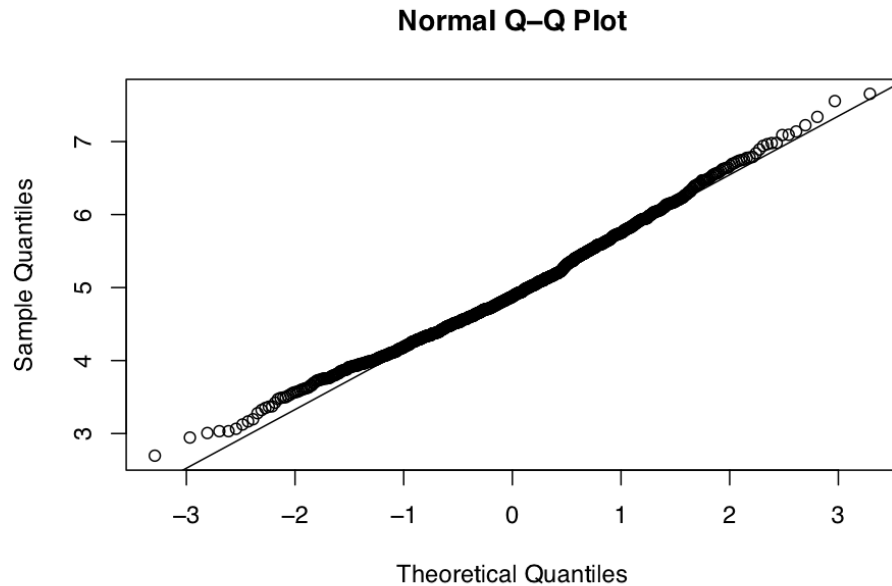
The theoretical variance is greater than the distribution's variance, so the distribution is less variable than its theoretical counterpart.

3. The distribution is approximately normal. The histogram of the distribution of sample means with a curve that describes the distribution as normal with the data's means and standard deviation.

Visually, the distribution seems normal.



To evaluate how close this distribution is to normal distribution, we can look at the sample means plotted in the quantile-quantile plot and against the line that represents normal distribution.



The samples means fall pretty much along the theoretical line, except at the quantiles farthest from the center on both sides, we see that the sample means deviate more from the line. Overall, we can say that the distribution is indeed approximately normal.

4. Evaluate the coverage of the confidence interval for $1/\lambda$: $\bar{X} \pm 1.96 \cdot (S/\sqrt{n})$. I calculated the 95% confidence intervals for each sample of 40 exponentials (1000 confidence intervals in total). This is what this table of intervals looks like:

```
## lower upper mean
## 1 3.035 6.264 4.650
## 2 3.266 6.596 4.931
## 3 3.603 8.273 5.938
## 4 3.337 8.746 6.041
## 5 3.544 6.397 4.970
## 6 2.868 5.804 4.336
```

I calculated how many of those confidence intervals cover the theoretical mean of $1/\lambda$.

```
## [1] 0.935
```

93.7% of the time, the confidence intervals of the data contain the true mean. But it is lower than the 95% confidence level used in the equation.

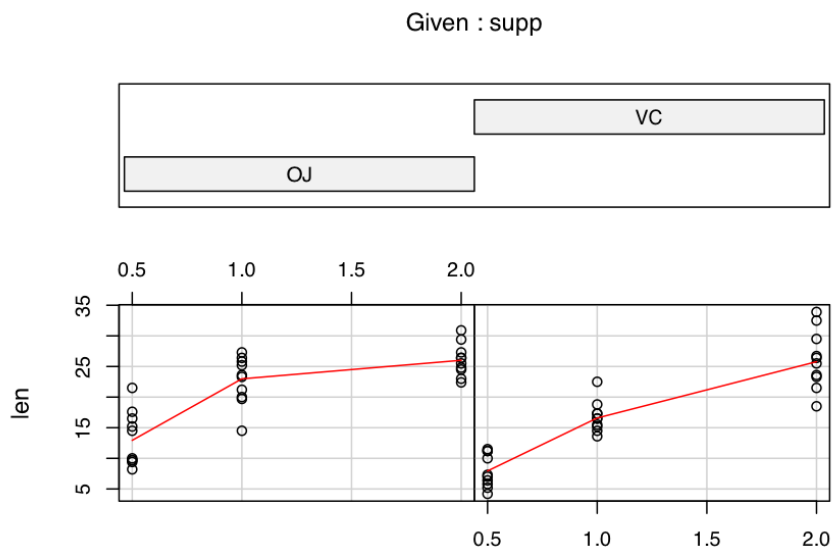
As we can see, this dataset of sample means has a distribution slightly different from the theoretical distribution, on the mean, variance, and confidence level. One guess is that the data is not large enough. If we increase the simulations to a number higher than 1000, we might get closer to the parameters of the true exponential distribution.

Part II. Analysis of the ToothGrowth Data

1. Load the ToothGrowth dataset, and
2. Provide a summary of the dataset.

```
##      len      supp      dose
## Min.   : 4.2    OJ:30    Min.   :0.50
## 1st Qu.:13.1    VC:30    1st Qu.:0.50
## Median :19.2                    Median :1.00
## Mean   :18.8                    Mean   :1.17
## 3rd Qu.:25.3                    3rd Qu.:2.00
## Max.   :33.9                    Max.   :2.00
```

A visual overview of the data, which shows the length of growth by supplement and dosage.



ToothGrowth data: length vs dose, given type of supplement

3. Use confidence intervals and hypothesis tests to compare tooth growth by supp and dose.

Since different guinea pigs receive different supplement and dosage, I assume that the variance is unequal and cannot be pooled so I should use an independent t-test for all comparison groups. The tables below compare the standard deviation of the supplement type groups and dosage groups. We can see that the variances are indeed not equal across comparison groups.

```
##  supp mean  sd  N  se
## 1   OJ 20.66 6.606 30 1.206
## 2   VC 16.96 8.266 30 1.509
```

```
##   dose mean    sd  N    se
## 1  0.5 10.61 4.500 20 1.0062
## 2  1.0 19.73 4.415 20 0.9873
## 3  2.0 26.10 3.774 20 0.8439
```

Comparing supplement types orange juice (OJ) v. ascorbic acid (VC)

```
##
## Welch Two Sample t-test
##
## data: len by supp
## t = 1.915, df = 55.31, p-value = 0.06063
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.171 7.571
## sample estimates:
## mean in group OJ mean in group VC
##           20.66           16.96
```

The result is not significant at the 95% confidence level and p-values is greater than 0.05. It means that there is no significant difference in tooth growth between using orange juice and ascorbic acid.

Comparing vitamin C dosage level of 0.5, 1, and 2mg.

```
##
## Pairwise comparisons using t tests with non-pooled SD
##
## data: ToothGrowth$len and ToothGrowth$dose
##
##   0.5    1
## 1 2.5e-07 -
## 2 1.3e-13 1.9e-05
##
## P value adjustment method: holm
```

The table shows comparisons among the three dosage levels. All of the p-values are less than 0.05, which means that tooth growth among all three groups with different dosage levels are significantly different.

We can conclude that the dosage level of vitamin C rather than the supplement type makes significant difference in subjects' tooth growth. Assumptions?