

# BANA 7031 - Probability Models

## Probabilistic analysis of Computer Prices

By

Syed Imad Husain (M12958531)

### *Abstract*

*This report discusses the steps involved in approximating a parametric distribution to real life data. We have used the Computers dataset in R and tried to determine the distribution for Computer Prices. Our conclusion shows us that the Prices follow a Lognormal distribution with  $\text{Meanlog} = 7.67$  and  $\text{Sdlog} = 0.257$ . We also tested if the prices of Computers with different types of RAM is the same using Permutation test and it seems that is not the case. Finally, we also discussed the Bayesian Posterior Distribution of the Prices towards the end of the report.*

### Contents

Introduction.....	2
Exploring Candidate Models .....	2
Bootstrap C&F .....	2
Theoretical vs Empirical graphs.....	2
AIC Criterion .....	5
Kolmogorov-Smirnov Simulation Test .....	5
Empirical Cumulative Distribution Function.....	6
Maximum Likelihood Parameter Estimation .....	6
Bootstrap for Standard Error MLE and Confidence Interval.....	6
Permutation Method to compare Means of two Groups.....	7
Bayesian Posterior Density for Distribution Mean .....	8
Appendix .....	8
R Code.....	9

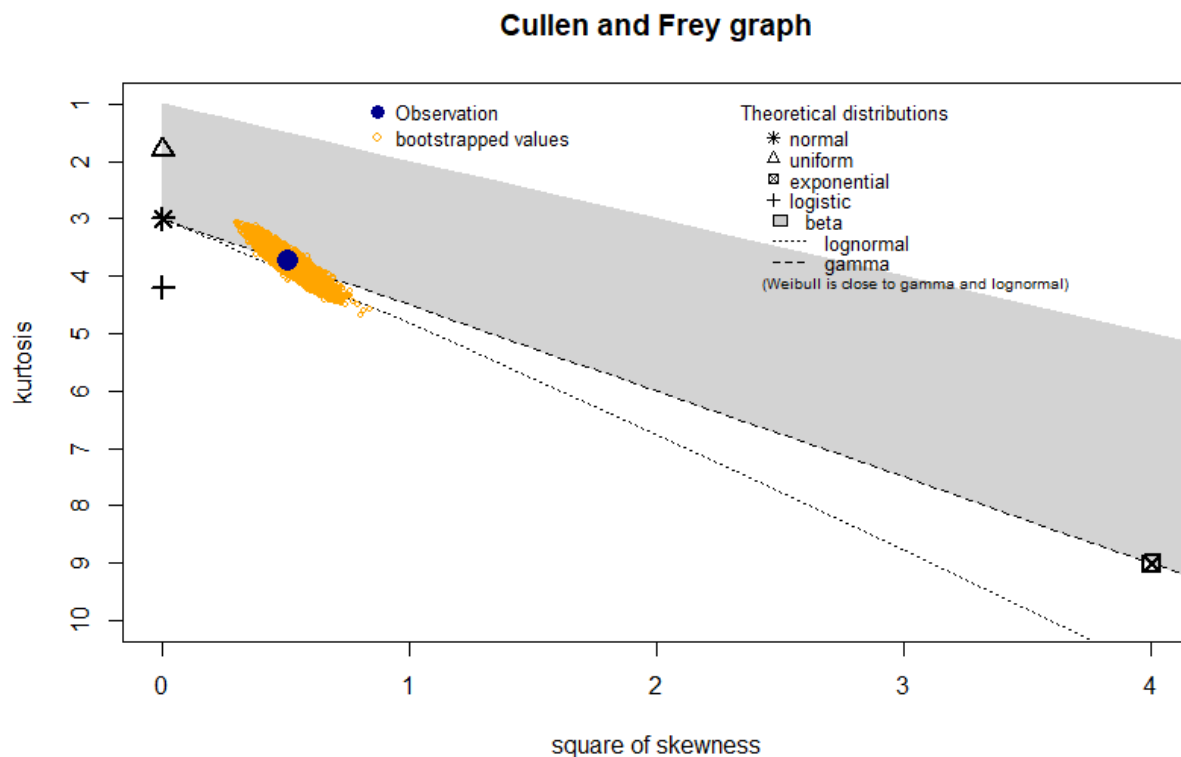
## Introduction

This artefact summarizes various types of analyses performed on Prices of Computers. The dataset being used is “Computers”<sup>[1]</sup>. The goal of this project is to come up with a probability distribution for computer prices. This is a parametric approach and all statistical work has been accomplished in R.

## Exploring Candidate Models

### Bootstrap C&F

There are several techniques for model selection. We will start our research by Moment-approach where we will try to evaluate the higher order (skewness<sup>2</sup>, kurtosis) space to narrow down on the set of probable distributions. The best way to visualize this analysis is the Cullen and Frey graph<sup>[2]</sup>(C&F Graph). It is only in our interest that we bootstrap the estimated kurtosis and skewness. Below C&F graph shows the result with 10,000 Bootstrapped values



The nearness of our observations(blue/yellow dots) to different theoretical distributions help us narrow the search space for candidate models. We conclude that Weibull, Normal, Lognormal and even Gamma are some of the probable distributions.

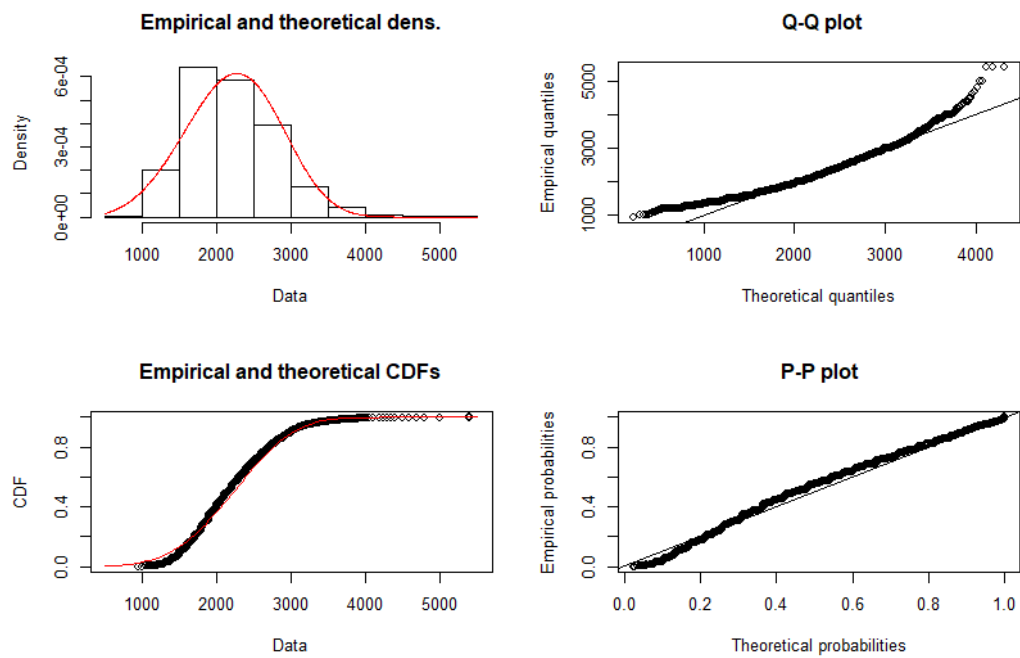
### Theoretical vs Empirical graphs

In order to further narrow our analysis, we will look at the following graphs comparing the “Theoretical vs Empirical” to visualize the fitting of the distribution to the data:

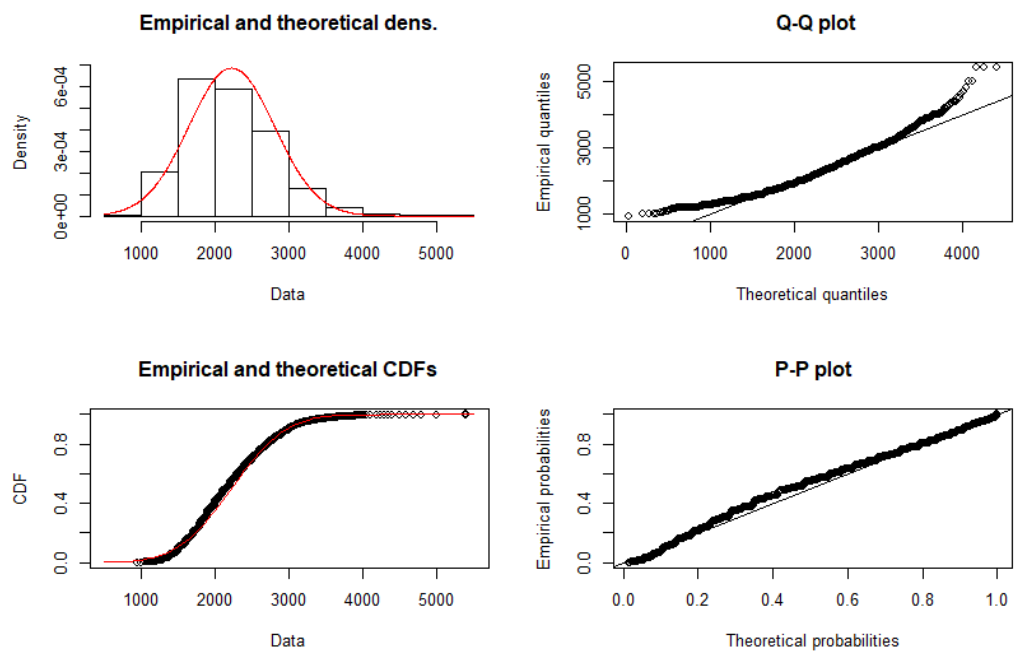
- Empirical PDF
- Empirical CDF
- Quantile Plots
- Percentile Plots

Following are the observations

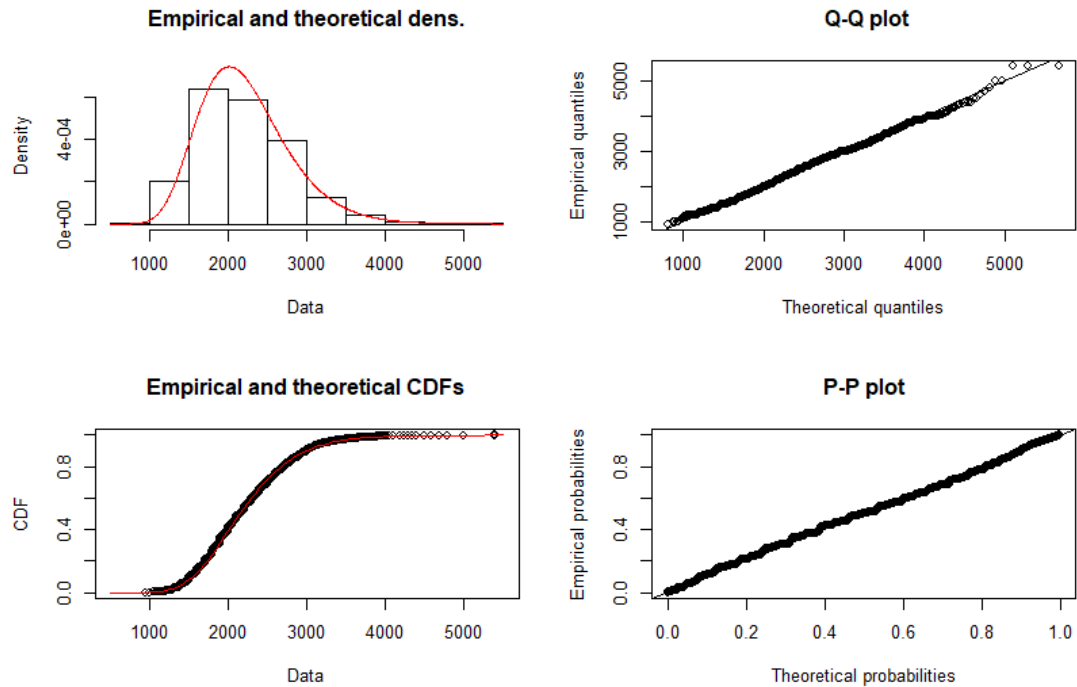
- Weibull Distribution



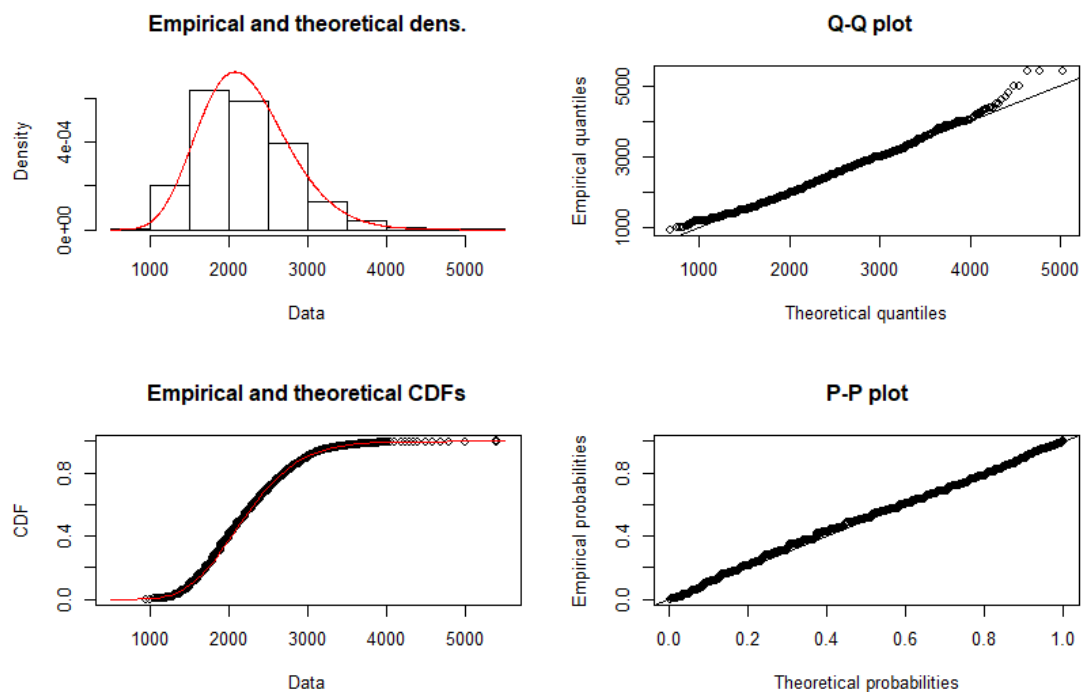
- Normal



- Lognormal



- Gamma



Based on these observations, we finalize **Log-Normal** to be the best candidate distribution. As we can see from above four graphs, the Q-Q plot and P-P plot, the empirical and theoretical dens. and Empirical and theoretical CDFs are the best among all four distribution.

## AIC Criterion

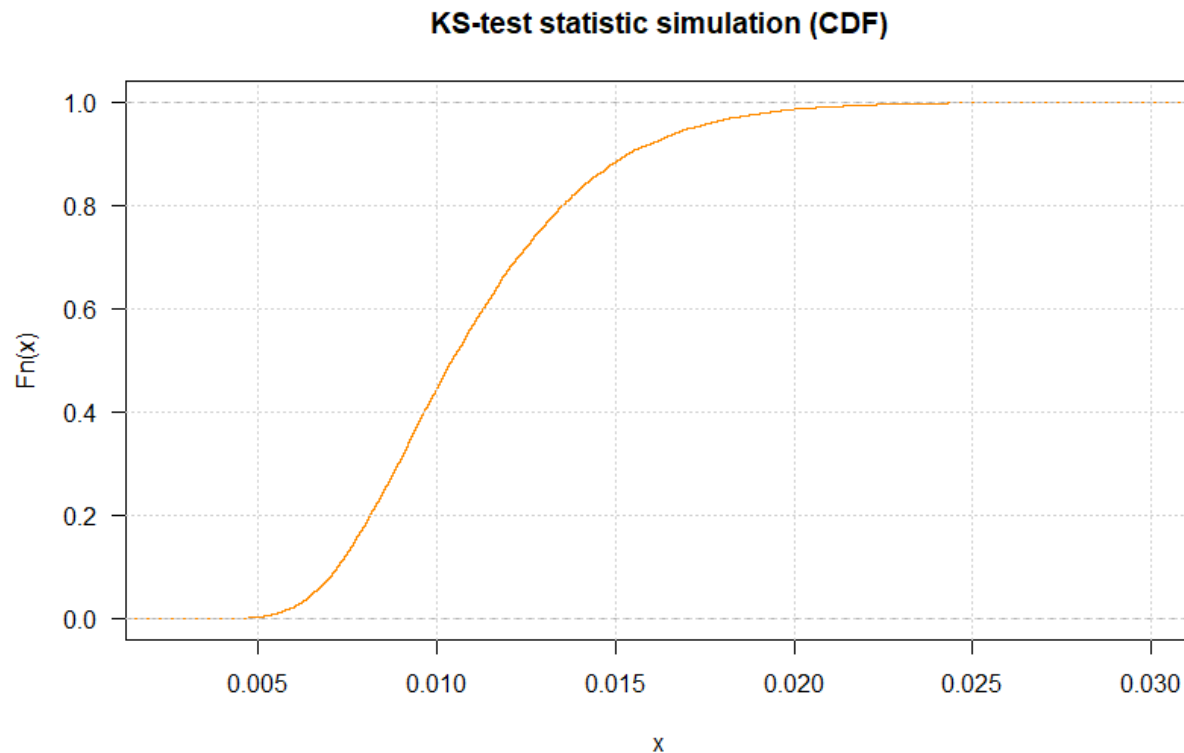
We will further use AIC to ascertain our assumption:

Distribution	AIC
Log-Normal	96841.68
Gamma	96914.93
Normal	97435
Weibull	97768.54

Distribution – Lognormal distribution with parameters Meanlog = 7.67177 Sdlog = 0.2579901.  
By using AIC, the smaller the better, the above table show us that the Log-Normal has the smallest AIC value.

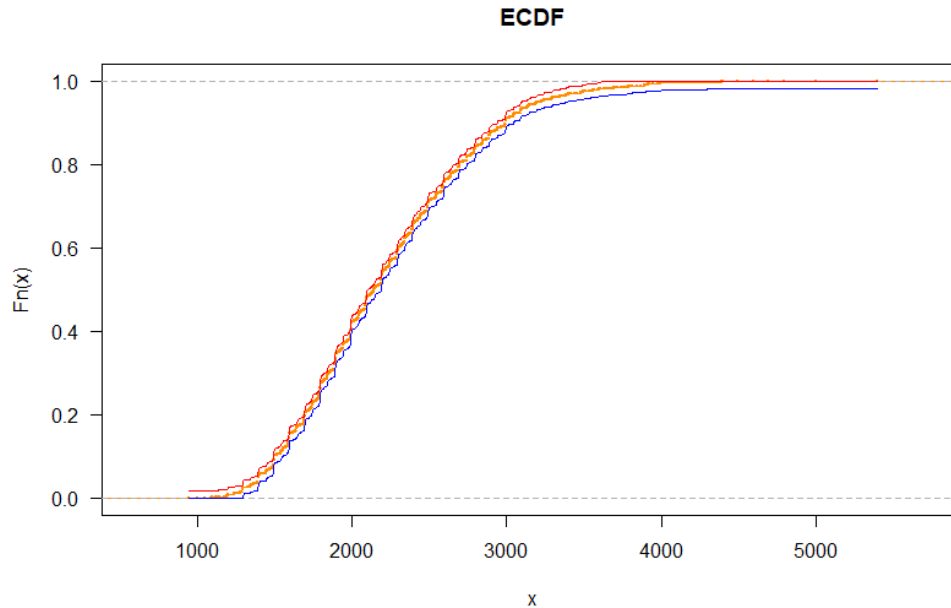
## Kolmogorov-Smirnov Simulation Test

After finalizing the distribution, we would like to check if this distribution is indeed a good fit. For this, we will simulate the KS statistic under the null-hypothesis that the parameters are 0. The ECDF for simulated KS statistic is plotted below and the p-value for the test is very less, hence we conclude that we have arrived at the correct distribution



## Empirical Cumulative Distribution Function

Since we now have a parametric distribution for the approximation, we will check its performance against an empirical CDF. Plot for the ECDF with 95% level of significance. The coverage of ECDF is 91.52%. This indicates that the ECDF confidence band was able to successfully contain the true value 91.52% times.



## Maximum Likelihood Parameter Estimation

The MLE for Lognormal distribution is given as follows, when calculated from data,

- Mean = 7.67177
- Standard Deviation = 0.2579901

$$\hat{\mu} = \frac{\sum_{i=1}^n \ln(X_i)}{n} \text{ and}$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \left( \ln(X_i) - \frac{\sum_{i=1}^n \ln(X_i)}{n} \right)^2}{n}.$$

## Bootstrap for Standard Error MLE and Confidence Interval

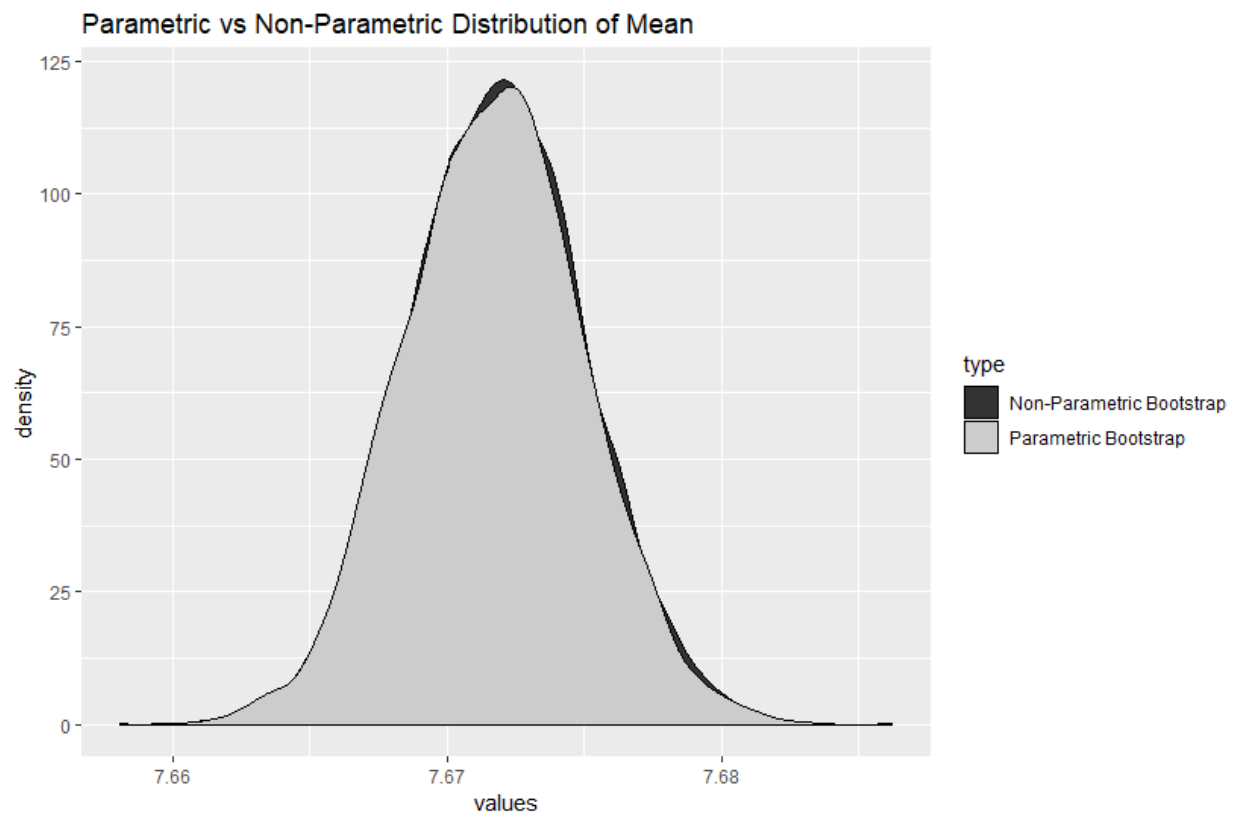
We perform Bootstrap to calculate the confidence interval and standard error for the MLE of the Mean. The Standard error from Bootstrap is

- Parametric = 0.003260
- Non-Parametric = 0.003284

We also calculated the following confidence intervals which all had a lower bound of 7.67 and an upper bound of 7.68 :

- Normal for Parametric
- Quantile for Non-Parametric
- Normal for Non-Parametric
- Pivotal for Non-Parametric

Below is the Bootstrap distribution from the two methods and they are very similar:



## Permutation Method to compare Means of two Groups

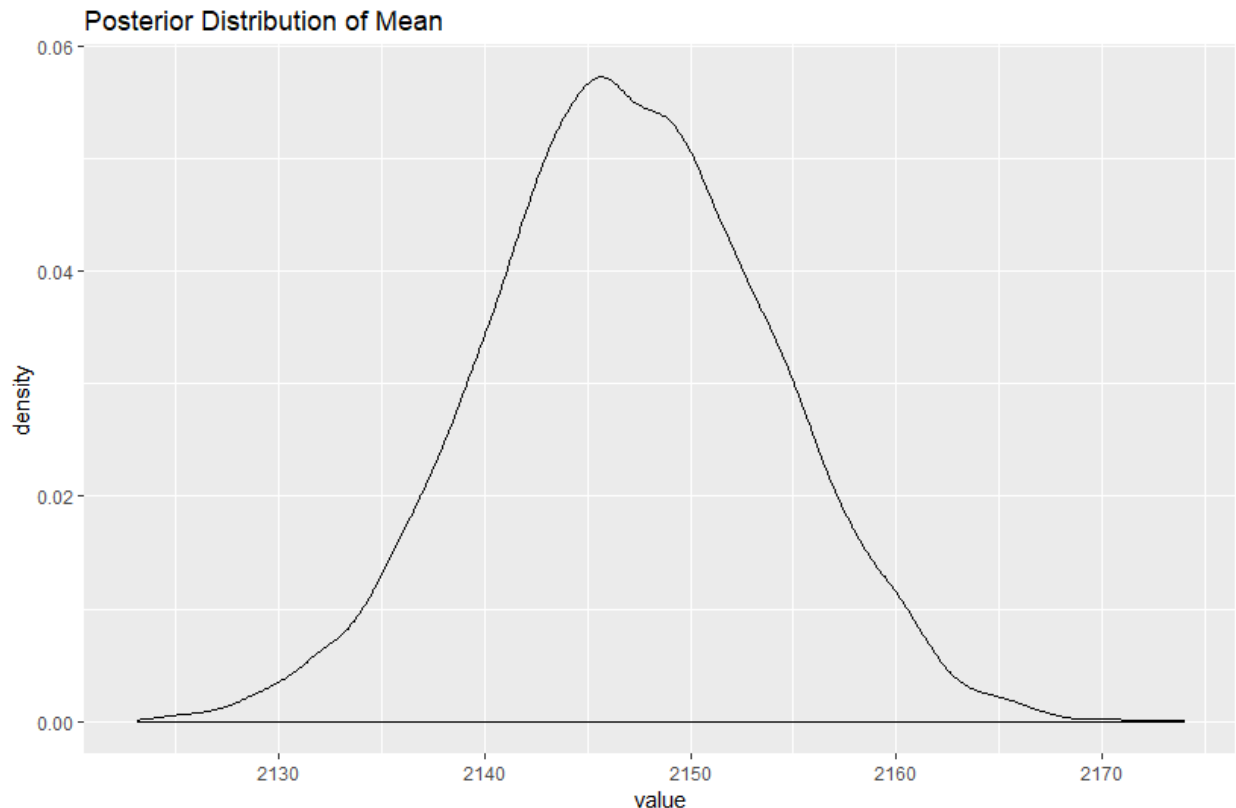
Here we try to check if the prices of computers with 4GB RAM are statistically different from 8GB RAM. In order to test this out, we perform the comparison through Permutation test.

- $H_0$  : Difference of Mean = 0
- $H_a$  : Difference of Mean  $\neq$  0

The p-value comes out to be very small and hence we reject the null hypothesis, thereby concluding that the prices of the two group are significantly different.

## Bayesian Posterior Density for Distribution Mean

There is no difference between observing a normal sample and a log-normal sample as one can be turned into the other by a log or exponential transform. We may hence use the ultra-classic Gaussian model framework. With a Prior for the Mean as  $f(\mu) = 1$  the Posterior Distribution is Log Normal  $(\bar{x}, \sigma^2/n)$ . Below is the simulation of Posterior distribution:



## Appendix

1. Computers Data set:
  - a. Description a cross-section from 1993 to 1995
  - b. number of observations: 6259
  - c. observation: goods
  - d. country: United States
  - e. Variables
    - i. **price** price in US dollars of 486 PCs
    - ii. **speed** clock speed in MHz
    - iii. **hd** size of hard drive in MB
    - iv. **ram** size of Ram in in MB
    - v. **screen** size of screen in inches
    - vi. **cd** is a CD-ROM present?
    - vii. **multi** is a multimedia kit (speakers, sound card) included?
    - viii. **premium** is the manufacturer was a "premium" firm (IBM, COMPAQ)?
    - ix. **ads** number of 486 price listings for each month



- x. **trend** time trend indicating month starting from January of 1993 to November of 1995.
  - f. Source Stengos, T. and E. Zacharias (2005) "Intertemporal pricing and price discrimination: a semiparametric hedonic analysis of the personal computer market", Journal of Applied Econometrics, forthcoming
  - g. References Journal of Applied Econometrics data archive:  
<http://qed.econ.queensu.ca/jae/>
2. Probabilistic Techniques in Exposure Assessment by **Cullen**, Alison C., **Frey**, H. Christopher

## R Code

```
3. #install.packages("fitdistrplus")
4. #install.packages("logspline")
5. library(fitdistrplus)
6. library(logspline)
7. library(Ecdat)
8. set.seed(007)
9.
10. x <- Computers$price
11.
12. descdist(x, discrete = FALSE, boot = 10000)
13.
14. fit.weibull <- fitdist(x, "weibull")
15. fit.norm <- fitdist(x, "norm")
16. fit.lognorm <- fitdist(x, "lnorm")
17. fit.gamma <- fitdist(x, "gamma")
18.
19. plot(fit.weibull)
20. plot(fit.norm)
21. plot(fit.lognorm)
22. plot(fit.gamma)
23.
24. fit.weibull$aic
25. fit.norm$aic
26. fit.lognorm$aic
27. fit.gamma$aic
28.
29. meanlog.x <- fit.lognorm$estimate[1]
30. sdlog.x <- fit.lognorm$estimate[2]
31. n <- length(x)
32. sim <- 10000
33. ks.sim <- rep(0, sim)
34.
35. for(i in 1:sim)
36. {
37.   resample.x <- rlnorm(n, meanlog = meanlog.x, sdlog = sdlog.x )
38.   ks.sim[i] <- as.numeric(ks.test(resample.x, "plnorm",
39.                                   meanlog = meanlog.x, sdlog = sdlog.x)$statistic)
40. }
41.
42. plot(ecdf(ks.sim), las = 1,
43.       main = "KS-test statistic simulation (CDF)", col = "darkorange", lwd =
44.         1.7)
45. grid()
```

```

45.
46. fit.final <- logspline(ks.sim)
47.
48. 1 - plogspline(ks.test(x,"plnorm"
49.                      ,meanlog = meanlog.x,sdlog = sdlog.x)$statistic
50.                      ,fit.final)
51.
52. x.ecdf <- ecdf(x)
53. plot(x.ecdf,
54.       las = 1,
55.       main = "ECDF",
56.       col = "darkorange",
57.       lwd = 1,cex=0.2
58.       )
59.
60. Alpha=0.05
61. Eps=sqrt(log(2/Alpha)/(2*n))
62. grid<-seq(min(x),max(x), length.out = 10000)
63. lines(grid, pmin(x.ecdf(grid)+Eps,1),col='red')
64. lines(grid, pmax(x.ecdf(grid)-Eps,0),col='blue')
65.
66. sum(plnorm(grid,meanlog.x,sdlog.x) >=pmax(x.ecdf(grid)-Eps,0) &
67. plnorm(grid,meanlog.x,sdlog.x) <=pmin(x.ecdf(grid)+Eps,1))/100
68.
69. mu.hat <- sum(log(x))/n
70. sd.hat <- sqrt(sum((log(x) - mu.hat)^2)/n)
71.
72. B <- 10000
73. mu.hat.star.p <- rep(0,B)
74. mu.hat.star.np <- rep(0,B)
75. for(i in 1:B)
76. {
77.   x.p <- rlnorm(n,mu.hat,sd.hat)
78.   x.np <- sample(x,size=n,replace=TRUE)
79.   mu.hat.star.p[i] <- sum(log(x.p))/n
80.   mu.hat.star.np[i] <- sum(log(x.np))/n
81. }
82. se.hat.boot.p <- sd(mu.hat.star.p)
83. se.hat.boot.np <- sd(mu.hat.star.np)
84.
85. par(mfrow <- c(2,1))
86. hist(mu.hat.star.p,main = "Parametric Method",prob=TRUE)
87. lines(density(mu.hat.star.p),col="red")
88. hist(mu.hat.star.np,main = "Parametric vs Non-Parametric Method",prob=TRUE)
89. lines(density(mu.hat.star.np),col="blue")
90.
91. library("tidyverse")
92. theta.np <- as.data.frame(mu.hat.star.np)
93. theta.p <- as.data.frame(mu.hat.star.p)
94.
95. names(theta.np) <- "values"
96. names(theta.p) <- "values"
97. theta.np$type <- 'Parametric Bootstrap'
98. theta.p$type <- 'Non-Parametric Bootstrap'
99. dis <- rbind(theta.p,theta.np)
100.   dis <- as.data.frame(dis)
101.   ggplot(dis, aes(values, fill = type)) + geom_density(alpha = 1) +

```

```

102.     scale_fill_grey() +
103.     ggtitle("Parametric vs Non-Parametric Distribution of Mean")
104.
105.     #normal ci, at 95% we get z.alpha/2 approx equal to 2
106.     normal.np<-c(mu.hat-2*se.hat.boot.np, mu.hat+2*se.hat.boot.np)
107.     #pivotal ci at 95%
108.     pivotal.np<-c(2*mu.hat-quantile(mu.hat.star.np,0.975),
109.                   2*mu.hat-quantile(mu.hat.star.np,0.025))
110.     #quantile ci at 95%
111.     quantile.np<-quantile(mu.hat.star.np, c(0.025, 0.975))
112.
113.     normal.p<-c(mu.hat-2*se.hat.boot.p, mu.hat+2*se.hat.boot.p)
114.
115.     y.4 <- Computers$price[Computers$ram=="4"]
116.     y.8 <- Computers$price[Computers$ram=="8"]
117.
118.     df <- c(y.4,y.8)
119.     length <- length(y.4)+length(y.8)
120.     options(expressions=1e5)
121.     B <- 10000
122.     perm.matrix <- replicate(B, sample(df,length,replace=F))
123.     perm.T <- rep(0,B)
124.     for(i in 1:B){perm.T[i] <- abs(mean(perm.matrix[1:length(y.4),i])-
125.     mean(perm.matrix[length(y.4):length(y.8),i]))}
126.     print(paste("The p-value for permutation test is",p.value <-
127.     mean(perm.T==0)))
128.
129.     sim <- 1000
130.     posterior.x <- as.data.frame(
131.         rlnorm(sim,meanlog = mu.hat,sdlog = sd.hat/sqrt(n)))
132.     names(posterior.x) <- "value"
133.     ggplot(posterior.x, aes(x=value,fill = value)) + geom_density(alpha = 1)
134.     +
135.     scale_fill_grey() +
136.     ggtitle("Posterior Distribution of Mean")

```