# Visualisation: Assignment 1

Students Name: XXXX XXXX Roll Number: MDS2021xx

Dead Line : 23 Nov 2021

## Instruction:

- Work on the 'Assignment 1.Rmd' file. Compile the file as pdf. Submit only the pdf file in moodle.
- If you want to do the work on Google colab, then please share the Colab link on the moodle.
- There are four problems.
- **Total 10 points**

## Problem 1 (3 points)

**Problem Statement:** Write an `R` function which will test Central Limit Theorem.

- Assume the underlying population distribution follow Poisson distribution with rate parameter $\lambda$
- We want to estimate the unknown $\lambda$ with the sample mean

$$\hat{\lambda} = \frac{1}{n}\sum_{i=1}^{n} X_i$$

- The exact sampling distribution of $\hat{\lambda}$ is unknown
- But CLT tells us that as sample size $n$ increases the sampling distribution of $\hat{\lambda}$ can be approximated by Gaussian distribution.

**Input in the function:** * n: sample size * $\lambda$ : rate parameter * N: simulation size

**Output from the function:**

- Histogram of the sampling distribution
- QQ-plot

**Test cases:** * case 1 a: $\lambda = 0.7$, n=10, N=5000 * case 1 b: $\lambda = 0.7$, n=30, N=5000 * case 1 c: $\lambda = 0.7$, n=100, N=5000 * case 1 c: $\lambda = 0.7$, n=300, N=5000

- case 2 a: $\lambda = 1.7$, n=10, N=5000
- case 2 b: $\lambda = 1.7$, n=30, N=5000
- case 2 c: $\lambda = 1.7$, n=100, N=5000
- case 2 c: $\lambda = 1.7$, n=300, N=5000

```
## write your R-function for problem 1 here
##
##
```

**Problem 2: (1 point)**

Consider the `JohnsonJohnson` dataset. The datset contains the Quarterly earnings (dollars) per Johnson & Johnson share 1960–80.

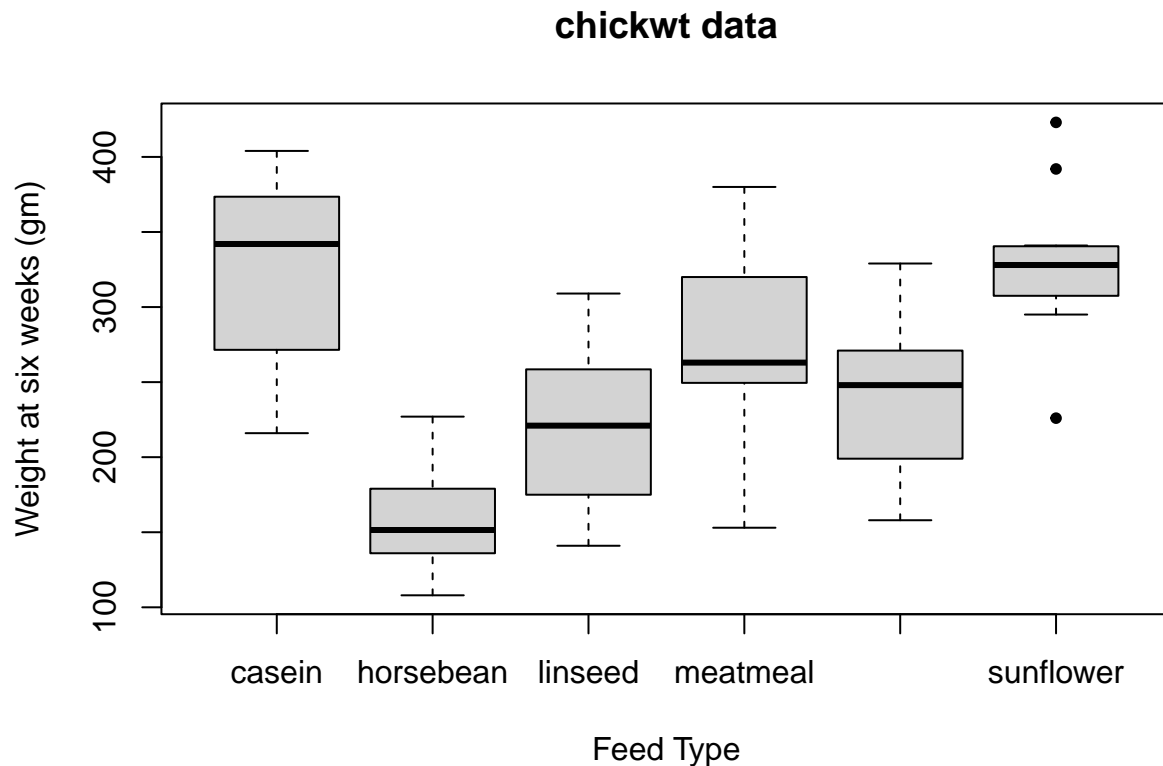    a) Draw the time series plot of Quarterly earnings in regular scale and log-scale using the `ggplot` (1 point)

```
head(JohnsonJohnson)
```

```
## [1] 0.71 0.63 0.85 0.44 0.61 0.69
```

**Problem 3: (2 points)**

- An experiment was conducted to measure and compare the effectiveness of various feed supplements on the growth rate of chickens.

- Following R-code is a standard side-by-side boxplot showing effect of feed supplements on the growth rate of chickens.

```
boxplot(weight~feed,data=chickwts,pch=20
        ,main = "chickwt data"
        ,ylab = "Weight at six weeks (gm)"
        ,xlab = "Feed Type")
```



**chickwt data**

    a) Reproduce the same plot using the `ggplot`; while fill each boxes with different colour. (1 point)

    b) In addition draw probability density plot for weights of chicken's growth by each feed seperately using the `ggplot`. Draw this plot seperately. (1 point)

**Problem 4: (4 points)**

- Consider the monthly data on the price of frozen orange juice concentrate in the orange-growing region of Florida.
- The data is available in `FrozenJuice` dataset of the `AER` package.
- We want to compare the average of price between decade of 1980's and 1990's. So we split the data into two

```
library(AER)
```

```
## Loading required package: car
```

```
## Loading required package: carData
```

```
## Loading required package: lmtest
```

```
## Loading required package: zoo
```

```
##
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
##
##     as.Date, as.Date.numeric
```

```
## Loading required package: sandwich
```

```
## Loading required package: survival
```
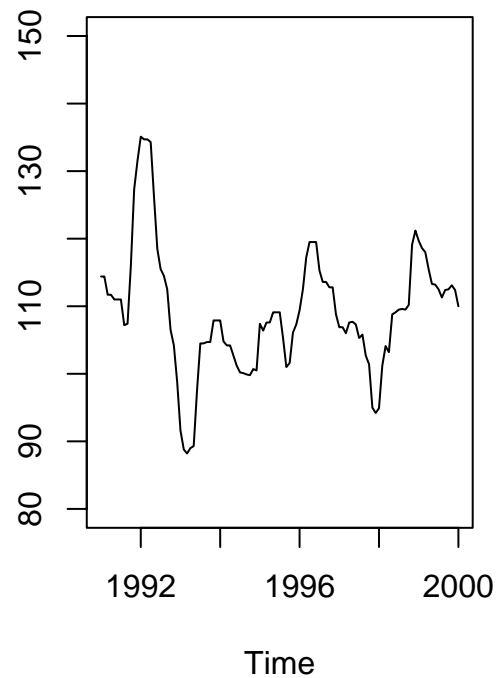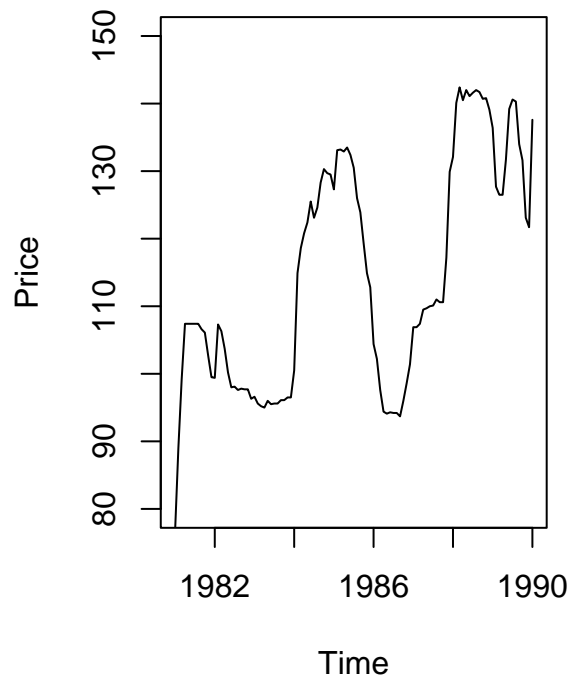
```
data("FrozenJuice")

data_80_90=window(FrozenJuice,start=1981,end=1990)
data_90_2K=window(FrozenJuice,start=1991,end=2000)
par(mfrow=c(1,2))
plot(data_80_90[,'price'],ylim=c(80,150),ylab='Price')
plot(data_90_2K[,'price'],ylim=c(80,150),ylab='')
```

- Generally it is believed that the price of the product increases over time due to inflation effect. So we expect that the average price during 1991-2000 would be higher than the 1981-1990.

The mean and standard deviation of price is estimates as

```
n1 = nrow(data_80_90)
cat('number of samples in 80s decade: ',n1,'\n')
```

```
## number of samples in 80s decade:  109
```

```
m1 = mean(data_80_90[,'price'])
s1 = sd(data_80_90[,'price'])
cat('mean and sd for 80s decade','\n')
```

```
## mean and sd for 80s decade
```

```
round(c(mean = m1,sd = s1),2)
```

```
##   mean     sd
## 114.32  16.88
```

```
n2 = nrow(data_90_2K)
cat('number of samples in 90s decade: ',n2,'\n')
```

```
## number of samples in 90s decade:  109
```

```
m2 = mean(data_90_2K[,'price'])
s2 = sd(data_90_2K[,'price'])
cat('mean and sd for 90s decade','\n')
```

```
## mean and sd for 90s decade
```

```
round(c(mean = m2,sd = s2),2)
```

```
##   mean     sd
## 109.14   9.25
```

```
round(c(mean = m2,sd = s2),2)
```

```
##   mean     sd
## 109.14   9.25
```

- The sample size for both decades are more than 100. So we can assume that CLT will kick-in.

a) If $\bar{X}_1$ and $\bar{X}_2$ are the sample mean of the price the two decades, plot the sampling distributions of sample mean for both decades on the same graph. (1 point)
b) Simulate the $\bar{X}_1$ and $\bar{X}_2$ from respective sampling distribution, then calculate the difference.

$$d = \bar{X}_1 - \bar{X}_2$$

   Simulate $d$; 5000 times. (1 point)
c) Calculate $P(d < 0)$ as

$$\hat{P}(d < 0) = \frac{\text{number of d<0}}{5000}$$

   and draw the histogram of $d$ and marked the area where $d < 0$ (1 point)
d) Based on the analysis, what is the chance that the average price of Juice for decade 1981-90 was same or less than the decade of 1991-2000? (1 point)