

# Probability and Statistics with R

## Assignment 2

Submission Nov 16-2022 (Wednesday)

**Note:** Below I explain how you collaborate on GitHub.

1. It will be group assignment.
2. A group would be of size at most 3. If you want to create a group size more than 3, you must take permission.
3. Decide among yourself and one of you create a GitHub repository for Probability Statistics Assignments.
4. In that repository add your group members as collaborator
5. Once you add your collaborator (or group members), create a folder and name it as **Assignment\_2**
6. In that folder you should have 2 folders **code** and **report**. And one **README.md** file. Write a brief report in **README.md** file.
7. For each problem, you should create a separate GitHub **issue**. All your discussion should be documented in the **issue**.
8. In the issue mention clearly, which group member is taking ownership of what problem?
9. The other member should **fork** the repository in their GitHub account.
10. Once you have your forked the main repository in your GitHub account - you should clone the repository in you local laptop or just download it as zip.
11. Once you develop the code - you should **commit** the code first in your repository and then **push** it.
12. Finally you make the **pull-request** in the final repository.
13. Once a member make a pull request, the other members have to review the code.
14. While reviewing the code the reviewer may have to download the code and run the code in his or her system and reproduce the result.
15. If the result is reproduced then she or he would accept and merge the code in final repository.
16. At the end you submit the link of the repository in the moodle.
17. The entire process will be evaluated.

### Problem 1

Suppose  $X$  denote the number of goals scored by home team in premier league. We can assume  $X$  is a random variable. Then we have to build the probability distribution to model the probability of number of goals. Since  $X$  takes value in  $\mathbb{N} = \{0, 1, 2, \dots\}$ , we can consider the geometric progression sequence as possible candidate model, i.e.,

$$S = \{a, ar, ar^2, ar^3, \dots\}.$$

But we have to be careful and put proper conditions in place and modify  $S$  in such a way so that it becomes proper probability distributions.

1. Figure out the necessary conditions and define the probability distribution model using  $S$ .
2. Check if mean and variance exists for the probability model.
3. Can you find the analytically expression of mean and variance.
4. From historical data we found the following summary statistics

mean	median	variance	total number of matches
1.5	1	2.25	380

Using the summary statistics and your newly defined probability distribution model find the following:

- a. What is the probability that home team will score at least one goal?
  - b. What is the probability that home team will score at least one goal but less than four goal?
5. Suppose on another thought you want to model it with off-the shelf Poisson probability models. Under the assumption that underlying distribution is Poisson probability find the above probabilities, i.e.,
- a. What is the probability that home team will score at least one goal?
  - b. What is the probability that home team will score at least one goal but less than four goal?
6. Which probability model you would prefer over another?
7. Write down the likelihood functions of your newly defined probability models and Poisson models. Clearly mention all the assumptions that you are making.

## Problem 2 : Simulation Study to Understand Sampling Distribution

**Part A** Suppose  $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} \text{Gamma}(\alpha, \sigma)$ , with pdf as

$$f(x|\alpha, \sigma) = \frac{1}{\sigma^\alpha \Gamma(\alpha)} e^{-x/\sigma} x^{\alpha-1}, \quad 0 < x < \infty,$$

The mean and variance are  $E(X) = \alpha\sigma$  and  $Var(X) = \alpha\sigma^2$ . Note that **shape** =  $\alpha$  and **scale** =  $\sigma$ .

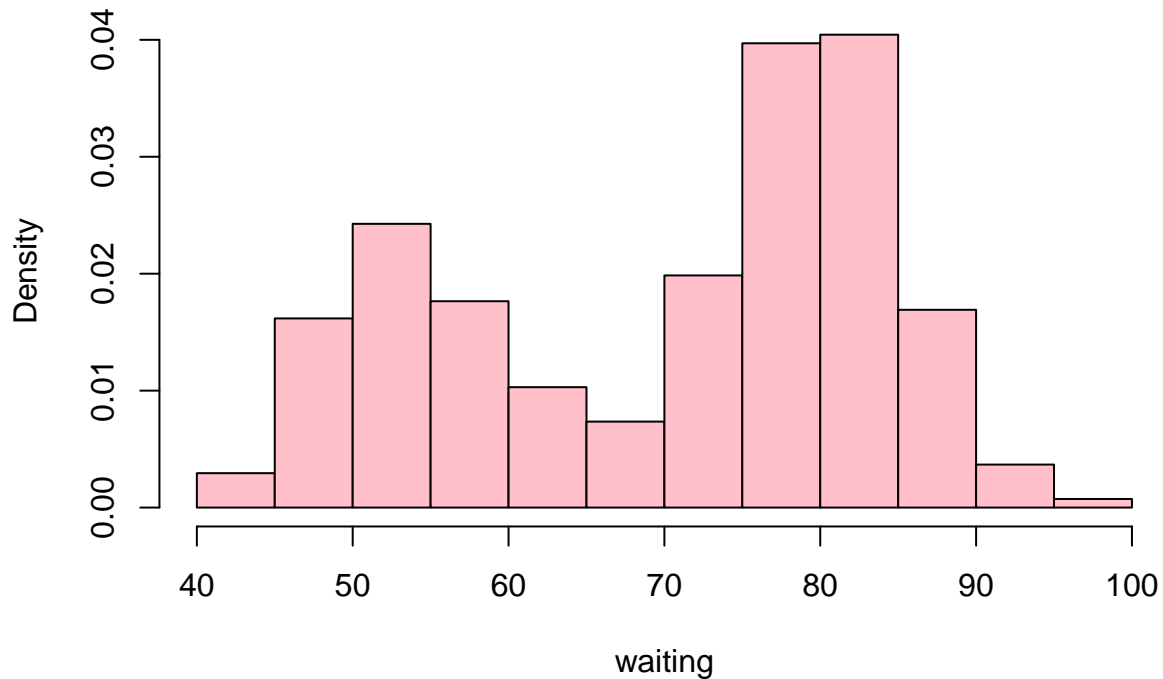
1. Write a **function** in R which will compute the MLE of  $\theta = \log(\alpha)$  using **optim** function in R. You can name it **MyMLE**
2. Choose **n=20**, and **alpha=1.5** and **sigma=2.2**
  - (i) Simulate  $\{X_1, X_2, \dots, X_n\}$  from **rgamma(n=20, shape=1.5, scale=2.2)**
  - (ii) Apply the **MyMLE** to estimate  $\theta$  and append the value in a vector
  - (iii) Repeat the step (i) and (ii) 1000 times
  - (iv) Draw histogram of the estimated MLEs of  $\theta$ .
  - (v) Draw a vertical line using **abline** function at the true value of  $\theta$ .
  - (vi) Use **quantile** function on estimated  $\theta$ 's to find the 2.5 and 97.5-percentile points.
3. Choose **n=40**, and **alpha=1.5** and repeat the (2).
4. Choose **n=100**, and **alpha=1.5** and repeat the (2).
5. Check if the gap between 2.5 and 97.5-percentile points are shrinking as sample size **n** is increasing?

*Hint:* Perhaps you should think of writing a single **function** where you will provide the values of **n**, **sim\_size**, **alpha** and **sigma**; and it will return the desired output.

### Problem 3: Analysis of faithful datasets.

Consider the faithful datasets:

```
attach(faithful)
hist(faithful$waiting, xlab = 'waiting', probability = T, col='pink', main='')
```



Fit following three models using MLE method and calculate **Akaike information criterion** (aka., AIC) for each fitted model. Based on AIC decides which model is the best model? Based on the best model calculate the following probability

$$\mathbb{P}(60 < \text{waiting} < 70)$$

(i) **Model 1:**

$$f(x) = p * \text{Gamma}(x|\alpha, \sigma_1) + (1 - p)N(x|\mu, \sigma_2^2), \quad 0 < p < 1$$

(ii) **Model 2:**

$$f(x) = p * \text{Gamma}(x|\alpha_1, \sigma_1) + (1 - p)\text{Gamma}(x|\alpha_2, \sigma_2), \quad 0 < p < 1$$

(iii) **Model 3:**

$$f(x) = p * \text{logNormal}(x|\mu_1, \sigma_1^2) + (1 - p)\text{logNormal}(x|\mu_1, \sigma_1^2), \quad 0 < p < 1$$

## Problem 4: Modelling Insurance Claims

Consider the **Insurance** datasets in the **MASS** package. The data given in data frame **Insurance** consist of the numbers of policyholders of an insurance company who were exposed to risk, and the numbers of car insurance claims made by those policyholders in the third quarter of 1973.

This data frame contains the following columns:

**District** (factor): district of residence of policyholder (1 to 4): 4 is major cities.

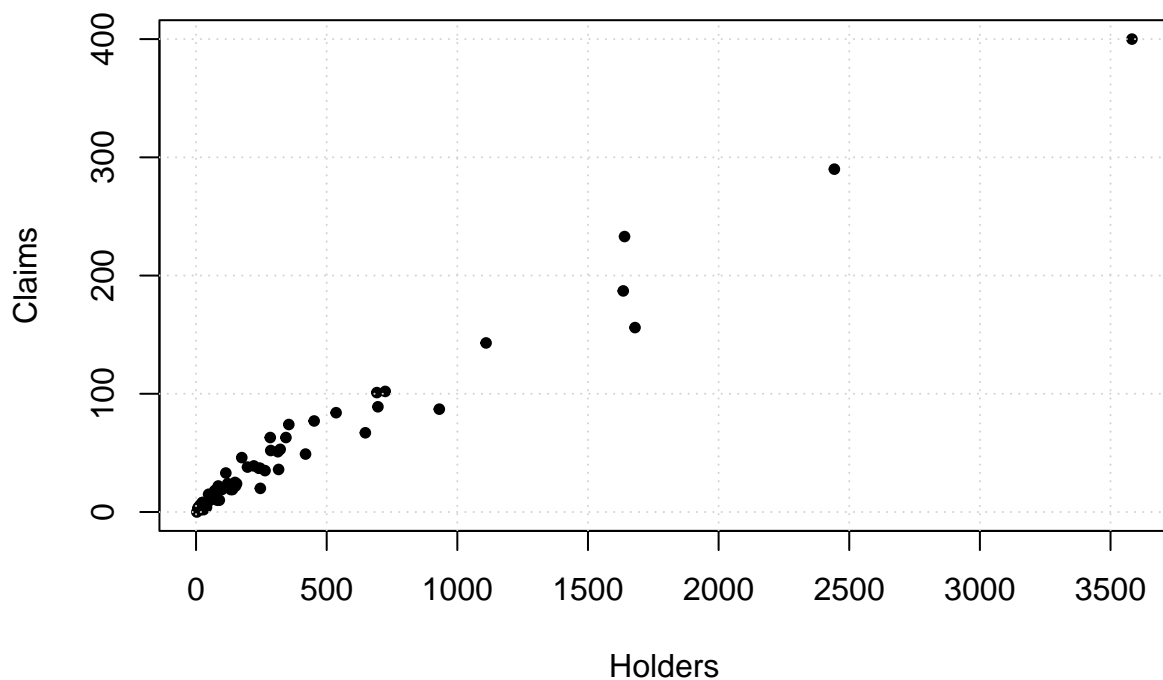
**Group** (an ordered factor): group of car with levels <1 litre, 1–1.5 litre, 1.5–2 litre, >2 litre.

**Age** (an ordered factor): the age of the insured in 4 groups labelled <25, 25–29, 30–35, >35.

**Holders** : numbers of policyholders.

**Claims** : numbers of claims

```
library(MASS)
plot(Insurance$Holders, Insurance$Claims
     ,xlab = 'Holders', ylab='Claims', pch=20)
grid()
```



**Note:** If you use built-in function like **lm** or any packages then no points will be awarded.

**Part A:** We want to predict the **Claims** as function of **Holders**. So we want to fit the following models:

$$\text{Claims}_i = \beta_0 + \beta_1 \text{Holders}_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

Assume :  $\varepsilon_i \sim N(0, \sigma^2)$ . Note that  $\beta_0, \beta_1 \in \mathbb{R}$  and  $\sigma \in \mathbb{R}^+$ .

The above model can also be re-expressed as,

$$\begin{aligned}\text{Claims}_i &\sim N(\mu_i, \sigma^2), \text{ where} \\ \mu_i &= \beta_0 + \beta_1 \text{Holders}_i + \varepsilon_i, \quad i = 1, 2, \dots, n\end{aligned}$$

- (i) Clearly write down the negative-log-likelihood function in R. Then use `optim` function to estimate MLE of  $\theta = (\beta_0, \beta_1, \sigma)$
- (ii) Calculate **Bayesian Information Criterion** (BIC) for the model.

**Part B:** Now we want to fit the same model with change in distribution:

$$\text{Claims}_i = \beta_0 + \beta_1 \text{Holders}_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

Assume :  $\varepsilon_i \sim \text{Laplace}(0, \sigma^2)$ . Note that  $\beta_0, \beta_1 \in \mathbb{R}$  and  $\sigma \in \mathbb{R}^+$ .

- (i) Clearly write down the negative-log-likelihood function in R. Then use `optim` function to estimate MLE of  $\theta = (\beta_0, \beta_1, \sigma)$
- (ii) Calculate **Bayesian Information Criterion** (BIC) for the model.

**Part C:** We want to fit the following models:

$$\begin{aligned}\text{Claims}_i &\sim \text{LogNormal}(\mu_i, \sigma^2), \text{ where} \\ \mu_i &= \beta_0 + \beta_1 \log(\text{Holders}_i), \quad i = 1, 2, \dots, n\end{aligned}$$

Note that  $\beta_0, \beta_1 \in \mathbb{R}$  and  $\sigma \in \mathbb{R}^+$ .

- (i) Clearly write down the negative-log-likelihood function in R. Then use `optim` function to estimate MLE of  $\theta = (\alpha, \beta, \sigma)$
- (ii) Calculate **Bayesian Information Criterion** (BIC) for the model.

**Part D:** We want to fit the following models:

$$\begin{aligned}\text{Claims}_i &\sim \text{Gamma}(\alpha_i, \sigma), \text{ where} \\ \log(\alpha_i) &= \beta_0 + \beta_1 \log(\text{Holders}_i), \quad i = 1, 2, \dots, n\end{aligned}$$

- (iii) Compare the BIC of all three models

## Problem 5: Computational Finance - Modelling Stock prices

Following piece of code download the prices of TCS since 2007

```
library(quantmod)

## Loading required package: xts
## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##   as.Date, as.Date.numeric
## Loading required package: TTR
## Registered S3 method overwritten by 'quantmod':
##   method           from
##   as.zoo.data.frame zoo
```

```
getSymbols('TCS.NS')
```

```
## Warning: TCS.NS contains missing values. Some functions will not work if objects  
## contain missing values in the middle of the series. Consider using na.omit(),  
## na.approx(), na.fill(), etc to remove or replace them.
```

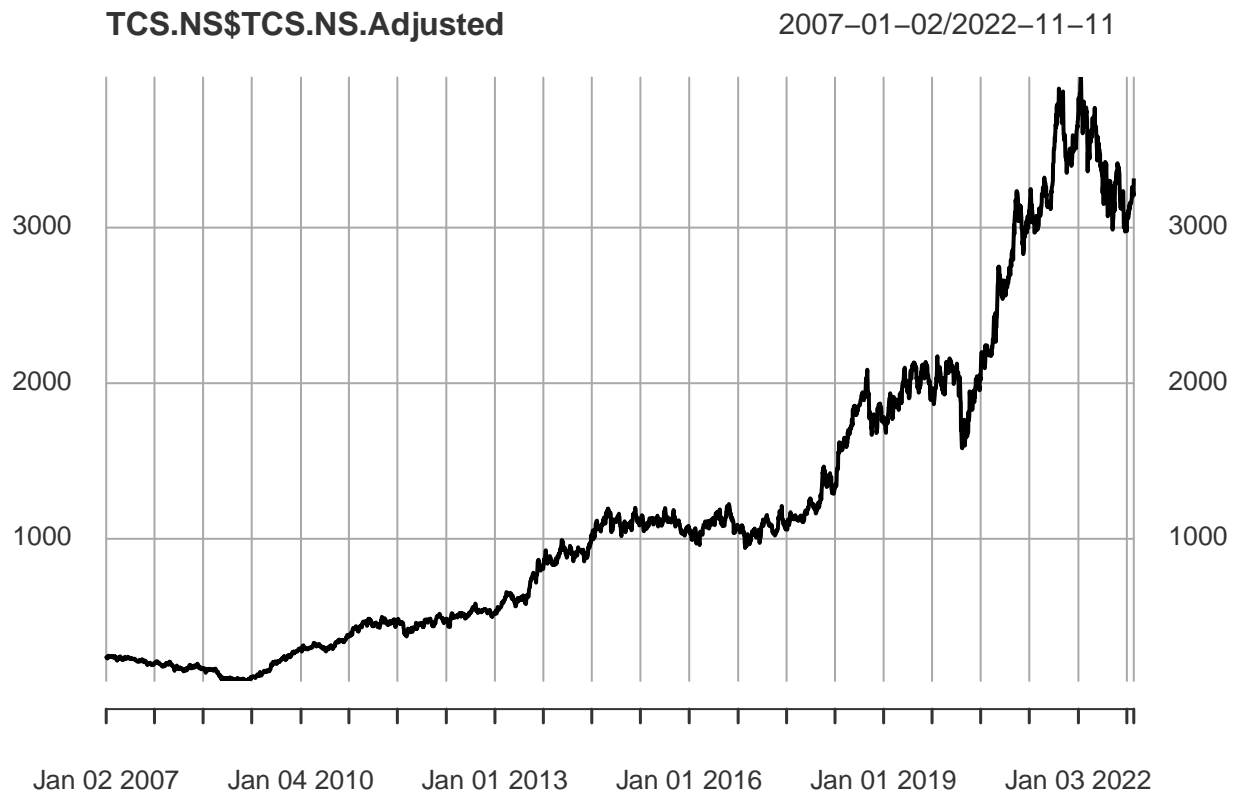
```
## [1] "TCS.NS"
```

```
tail(TCS.NS)
```

```
##           TCS.NS.Open TCS.NS.High TCS.NS.Low TCS.NS.Close TCS.NS.Volume  
## 2022-11-03      3228.05      3228.05      3195.00      3206.75      1422652  
## 2022-11-04      3217.00      3220.05      3166.15      3217.40      1464013  
## 2022-11-07      3229.00      3242.80      3195.10      3233.70      1474498  
## 2022-11-09      3249.80      3249.80      3201.65      3216.05      1162267  
## 2022-11-10      3170.00      3225.00      3170.00      3205.65      1573092  
## 2022-11-11      3269.60      3341.60      3255.05      3315.95      3265394  
##           TCS.NS.Adjusted  
## 2022-11-03           3206.75  
## 2022-11-04           3217.40  
## 2022-11-07           3233.70  
## 2022-11-09           3216.05  
## 2022-11-10           3205.65  
## 2022-11-11           3315.95
```

Plot the adjusted close prices of TCS

```
plot(TCS.NS$TCS.NS.Adjusted)
```



**Download the data of market index Nifty50.** The Nifty 50 index indicates how the over all market has done over the similar period.

```
getSymbols('^NSEI')
```

```
## Warning: ^NSEI contains missing values. Some functions will not work if objects
## contain missing values in the middle of the series. Consider using na.omit(),
## na.approx(), na.fill(), etc to remove or replace them.
```

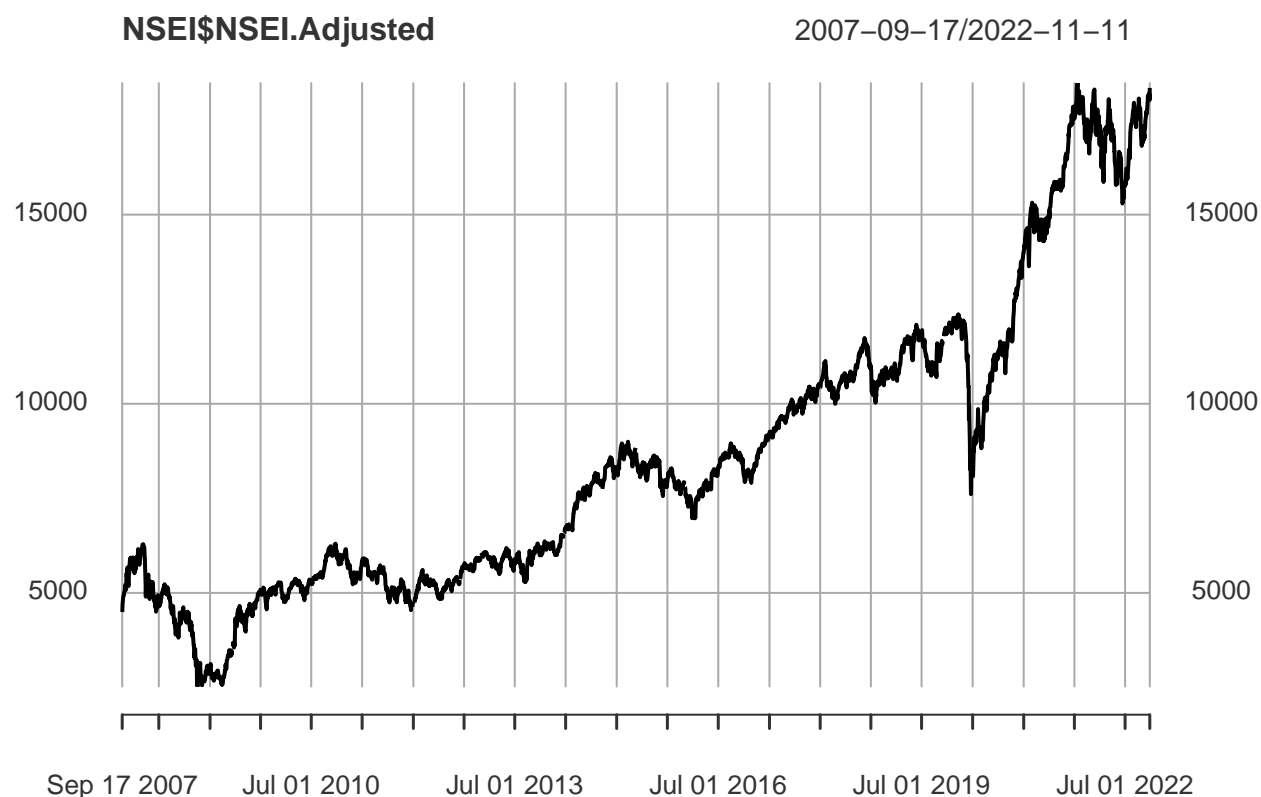
```
## [1] "^NSEI"
```

```
tail(NSEI)
```

```
##           NSEI.Open NSEI.High NSEI.Low NSEI.Close NSEI.Volume NSEI.Adjusted
## 2022-11-03  17968.35  18106.3 17959.20  18052.70      213000      18052.70
## 2022-11-04  18053.40  18135.1 18017.15  18117.15      267900      18117.15
## 2022-11-07  18211.75  18255.5 18064.75  18202.80      314800      18202.80
## 2022-11-09  18288.25  18296.4 18117.50  18157.00      307200      18157.00
## 2022-11-10  18044.35  18103.1 17969.40  18028.20      256500      18028.20
## 2022-11-11  18272.35  18362.3 18259.35  18349.70      378500      18349.70
```

Plot the adjusted close value of Nifty50

```
plot(NSEI$NSEI.Adjusted)
```



## Log-Return

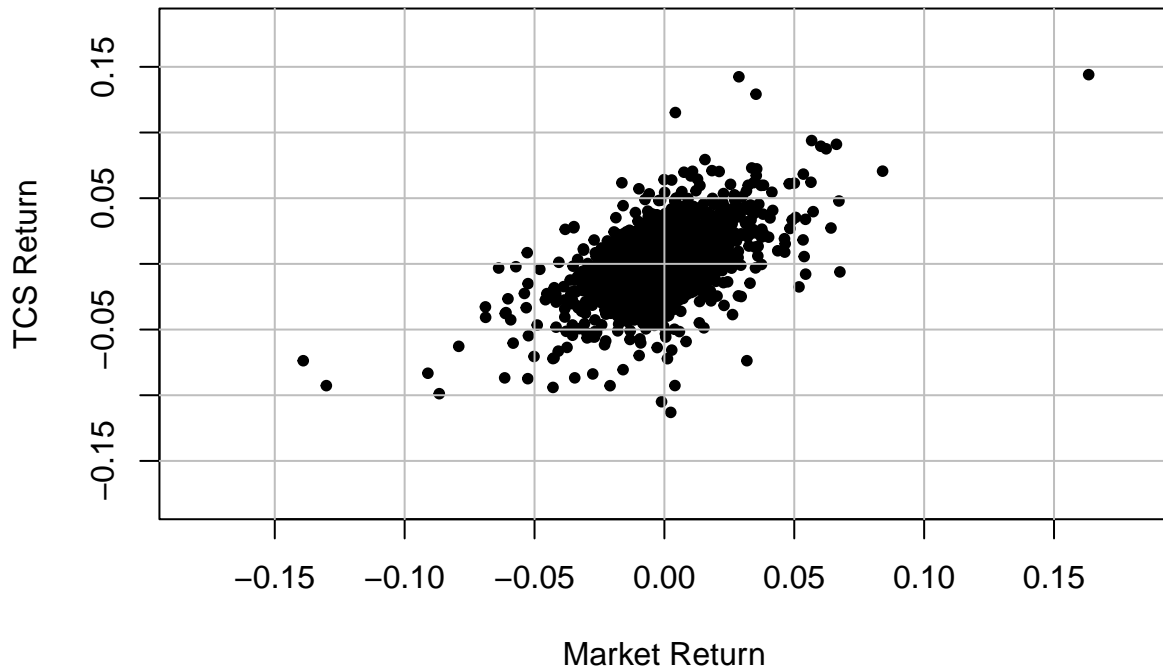
We calculate the daily log-return, where log-return is defined as

$$r_t = \log(P_t) - \log(P_{t-1}) = \Delta \log(P_t),$$

where  $P_t$  is the closing price of the stock on  $t^{th}$  day.

```
TCS_rt = diff(log(TCS.NS$TCS.NS.Adjusted))
Nifty_rt = diff(log(NSEI$NSEI.Adjusted))
retrn = cbind.xts(TCS_rt,Nifty_rt)
retrn = na.omit(data.frame(retrn))

plot(retrn$NSEI.Adjusted,retrn$TCS.NS.Adjusted
      ,pch=20
      ,xlab='Market Return'
      ,ylab='TCS Return'
      ,xlim=c(-0.18,0.18)
      ,ylim=c(-0.18,0.18))
grid(col='grey',lty=1)
```



- Consider the following model:

$$r_t^{TCS} = \alpha + \beta r_t^{Nifty} + \varepsilon,$$

where  $\mathbb{E}(\varepsilon) = 0$  and  $\text{Var}(\varepsilon) = \sigma^2$ .

1. Estimate the parameters of the models  $\theta = (\alpha, \beta, \sigma)$  using the method of moments type plug-in estimator discussed in the class.
2. Estimate the parameters using the `lm` built-in function of R. Note that `lm` using the OLS method.
3. Fill-up the following table



Parameters	Method of Moments	OLS
$\alpha$		
$\beta$		
$\sigma$		

4. If the current value of Nifty is 18000 and it goes up to 18200. The current value of TCS is Rs. 3200/-. How much you can expect TCS price to go up?