# Experiment 3.3

**Student Nam** : Mayank Kumar          **UI**  :20BCS1353
**Branch**: BE-CSE                      **Section/Group**: 20BCS_DM_705 A
**Semester**: 6th                       **Date of Performance**:
**Subject Name**: Data Mining Lab       **Subject Code**: 20CSP-376

**Aim**: Study of Outlier detection using R programming.

**Objective:** Data points far from the dataset's other points are considered outliers. This refers to the data values dispersed among other data values and upsetting the dataset's general distribution.

Effects of an outlier on model:

- The format of the data appears to be skewed.
- Modifies the mean, variance, and other statistical characteristics of the data's overall distribution.
- Leads to the model's accuracy level being biased.

**Script and Output**:

#creating the data containing 500 random values

data <- rnorm(500)

print(data)

#adding 10 random outliers to this data.

data[1:10] <- c(46,9,15,-90,42,50,-82,74,61,-32)

#draw boxpolot and an outlier is defined as a data point that is located outside the whiskers of the box plot.

boxplot(data)

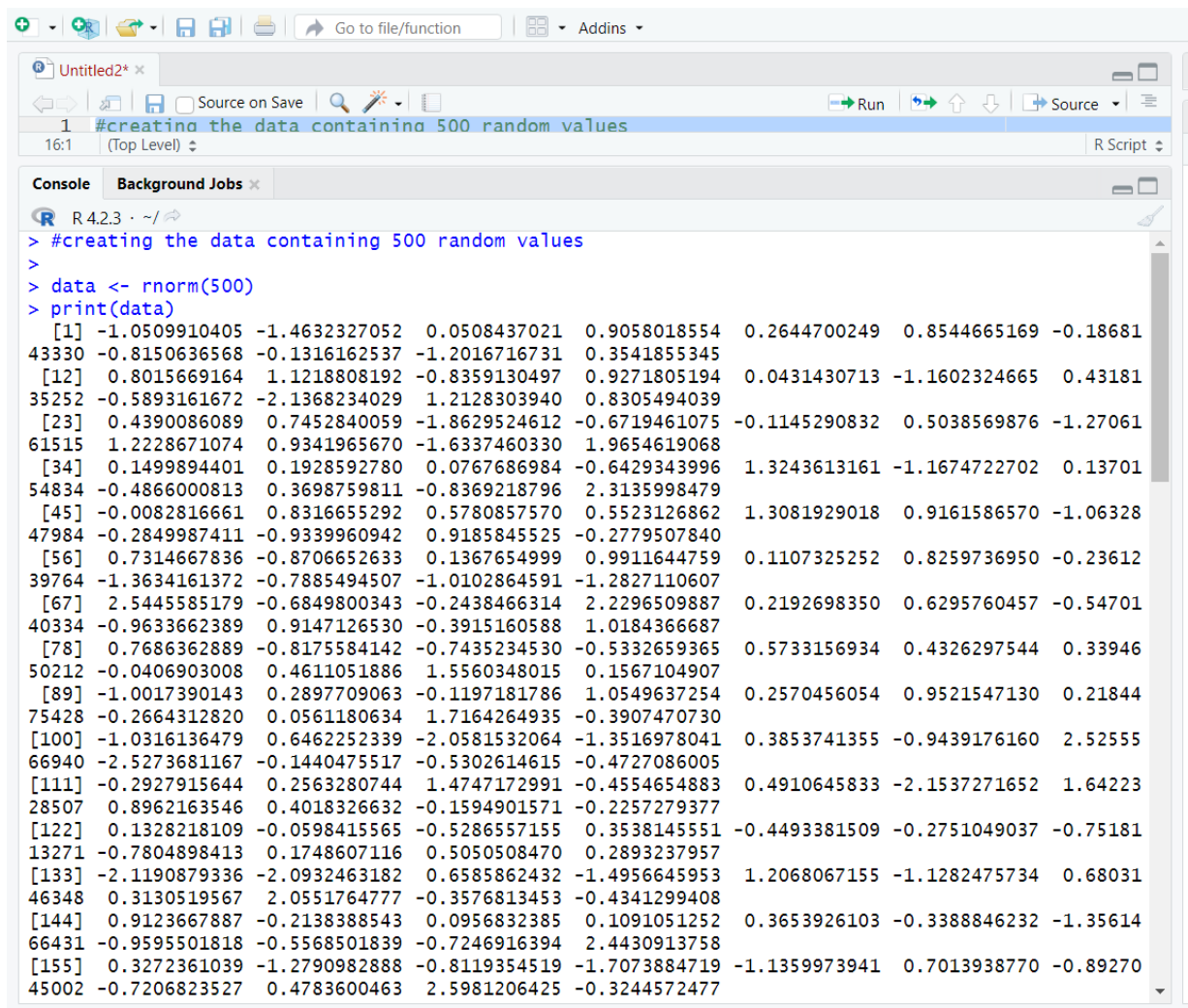#remove the outlier of the provided data boxplot.stats() function in R

data <- data[!data %in% boxplot.stats(data)$out]

#draw boxplot to verify whether ouliers removed or not

boxplot(data)

**OUTPUT:**

```
Untitled2* ×

   Source on Save                                              Run      Source
   1  #creating the data containing 500 random values
16:1   (Top Level)                                                      R Script

Console   Background Jobs ×

R 4.2.3 · ~/

[386]  0.3184278108   0.6077614987  -0.1399442884   0.2856267390  -0.2828708543  -0.8604843525   2.09029
86427  1.6797251307   1.4696083021  -0.7054707289  -0.1536420525
[397]  0.4201056849   0.5017141181   0.2160692129   0.6076228407   0.1047022229  -0.2093737402  -0.27079
25271  0.5565657225   2.1134504634  -1.5254271332   1.2507200372
[408] -0.5627236787   1.5676659970   1.5065637383   0.2715803427  -0.2328211332  -1.8212790271  -0.48522
80737 -0.3984773293   0.1819751019   0.6066645750  -1.0527947466
[419]  1.0211060598   0.8010015696  -0.4462740308   0.2552056396   0.3550136778   0.1499400366  -1.39384
47669  1.0874928242   1.4219486380   0.6032360547  -0.4741819261
[430] -0.4830995468   1.4433340416  -0.7566033587  -0.4598533764  -0.6581690589   0.5646506768  -0.64123
41256 -0.3294629345  -0.1211665393   0.0560386129   0.6342303681
[441] -0.9739237314   1.2669686857   1.8735877915   0.0584691661  -0.6766369722   0.0810428930  -1.28513
96595 -0.1296978675  -2.6075892168  -0.5358130440   0.6566164421
[452] -1.4000642070   0.2395365275  -1.4879477884  -1.3628601751  -0.3764113678  -2.0958710363  -0.81897
14866 -0.7854918242  -0.7183564473   0.1334205057  -0.0033126916
[463] -0.2819092781   1.8412405776   0.0554882345   0.1375625048  -1.1099579268   0.1214358996   0.38994
98093 -1.1310283595  -1.2565997451  -0.9407353915   0.3906926111
[474]  2.7670937860   0.8555548767  -0.6523407955  -0.3260157766  -0.0481155498  -0.6742881165  -0.67210
82528  1.2411610688  -0.7506700996   1.0306432936   0.0617763375
[485]  0.7611875336   1.5870086842   0.0525892434   0.6334560324  -0.8163511674  -0.2812928992   0.32488
16418 -0.2821754260  -0.5093303665  -0.1890192222  -1.7913864289
[496]  0.7497021854  -1.0208487285   0.0198523693  -0.9874702390  -0.6581323895
> #adding 10 random outliers to this data.
> data[1:10] <- c(46,9,15,-90,42,50,-82,74,61,-32)
>
> #draw boxpolot and an outlier is defined as a data point that is located outside the whiskers of
the box plot.
> boxplot(data)
>
> #remove the outlier of the provided data boxplot.stats() function in R
> data <- data[!data %in% boxplot.stats(data)$out]
>
> #draw boxplot to verify whether ouliers removed or not
> boxplot(data)
> |
```

**Learning Outcome:**

1. Learnt about Regression Analysis using R Programming.
2. Learnt about Simple Linear Regression.
3. Learnt about Multiple Linear Regression.