

# Líkindareikningur og Tölfræði STÆ203G

## Tölvuverkefni 3

Egill Ian Guðmundsson, 260693-2639

### Verkefni 1

Sólarhringsúrcoma hvers dags er mæld kl 9:00 á morgnana. Mælingar á mestu sólarhringsúrkoma (mm) á Fagurhólsmýri innan hvers árs fyrir árin 1924 til 2007 ( $n = 84$ ) er að finna í gagnaskránni `precipitation_fagurhm.txt`. Notið þessi gögn til að reikna stærðirnar sem eru gefnar fyrir neðan.

**1.a.** Reiknið meðaltal, miðgildi, dreifni, staðalfrávik, fyrsta fjórðungsmark, þriðja fjórðungsmark og fjórðungsbil fyrir úrtakið og setjið saman í eina töflu.

**Svar:** Við keyrum `precipitationStats.r` forritið sem fylgir með í viðauka til að fá eftirfarandi niðurstöðu:

```

Console C:/Users/egill/OneDrive/4Misseri/LikindiOgTolfraedi/R_verkefni/R_skjol/
> precipitationStats()
      Mean Median Variance Standard Deviation First Quartile Third Quartile Interquartile Range
value 71.72024   69.7 355.3445      18.85058         56.775         82.525          25.75
Mean value is larger than median value.
Range with 50 percent of values with median as mid-point: [56.775 ; 82.525]
      Mean Median Variance Standard Deviation First Quartile Third Quartile Interquartile Range
Logarithmic values 4.239824 4.24415 0.06630658      0.2575006         4.039089         4.413098          0.374009
> |

```

Þar sem öll gildin sem beðið var um koma fyrir.

**1.b.** Hvort er meðaltalið eða miðgildið stærra? Hvers vegna?

**Svar:** Eins og sést af töflunni að ofan er meðaltalið stærra en miðgildið. Þetta er væntanlega sökum þess að að eru nokkrir "útlagar" sem eru mun hærri en miðgildið en engir "útlagar" sem eru fyrir neðan miðgildið.

Áhrifin sem þetta hefur er að útlagarnir fyrir ofan hífa upp meðaltalið og skekkja það svo það verður hærri en miðgildið. Eins myndi miðgildið vera hærri en meðaltalið ef útlagarnir fyrir neðan miðgildið væru fleiri en útlagarnir fyrir ofan.

**1.c.** Á hvaða bili liggja 50% af mælingum sitt hvorum megin við miðgildið?

**Svar:** Svarið er einfaldlega bilið sem er milli fyrsta fjórðungsmarks og þriðja fjórðungsmarks. Eins og sést á mynd að ofan reiknast bilið [56.775 ; 82.525]

**1.d.** Reiknið logrann af hverju staki í gögnunum og reiknið sömu stærðir og hér að ofan í 1.a. en setjið gögnin í nýja töflu.

**Svar:** Aftur má athuga myndina til að sjá töfluna með gögnunum fyrir logra-gögnin:

```

Console C:/Users/egill/OneDrive/4Misseri/LikindiOgTolfraedi/R_verkefni/R_skjol/
> precipitationStats()
      Mean Median Variance Standard Deviation First Quartile Third Quartile Interquartile Range
value 71.72024   69.7 355.3445      18.85058         56.775         82.525          25.75
Mean value is larger than median value.
Range with 50 percent of values with median as mid-point: [56.775 ; 82.525]
      Mean Median Variance Standard Deviation First Quartile Third Quartile Interquartile Range
Logarithmic values 4.239824 4.24415 0.06630658      0.2575006         4.039089         4.413098          0.374009
> |

```

## Verkefni 2

Hér á að teikna myndir af gögnunum. Byggt á myndunum á að meta hvort að normaldreifingin lýsi gögnunum nægjanlega vel. Einnig á að meta hvort að normaldreifingin lýsi logranum af gögnunum nægjanlega vel. Teiknið eftirfarandi:

**2.a.** Teiknið árlega hámarkssólarhringsúrkomu á móti tíma. Á hvaða ári rigndi mest?

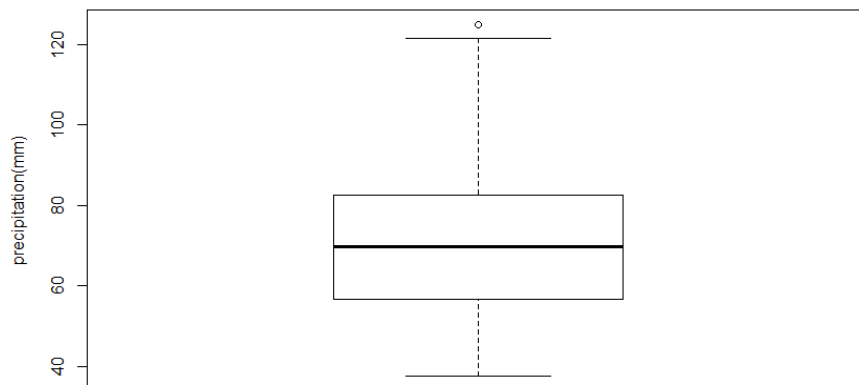
**Svar:** Við keyrum fall úr `precipitationGraphs.r` með skipun `precipitationGraphs(0,FALSE)` til að fá eftirfarandi mynd:



Forritið gefur síðan eftirfarandi úttak á skipanalínu:  
Maximum precipitation was 125 in the year 1936

**2.b.** Kassarit. Sýnir kassaritið einhverja útlaga?

**Svar:** Við keyrum sama fall og áðan með skipun `precipitationGraphs(1,FALSE)` til að fá eftirfarandi mynd:

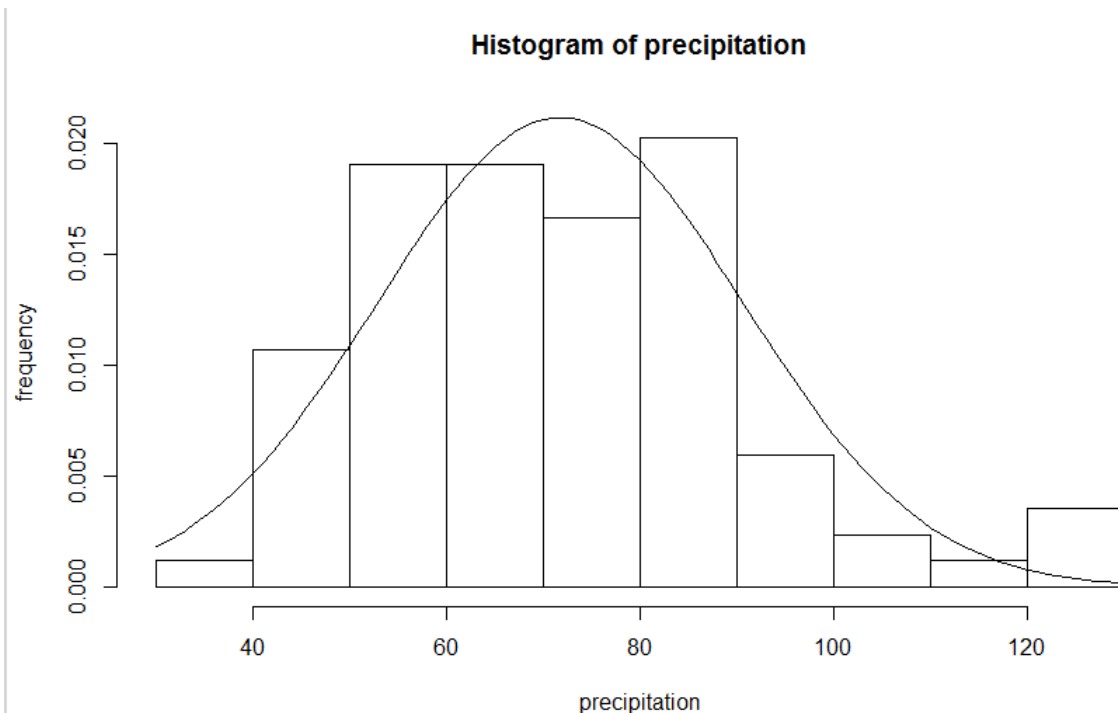


Einnig gefur forritið eftirfarandi úttak á skipanalínuna:

Number of outliers in boxplot: 2

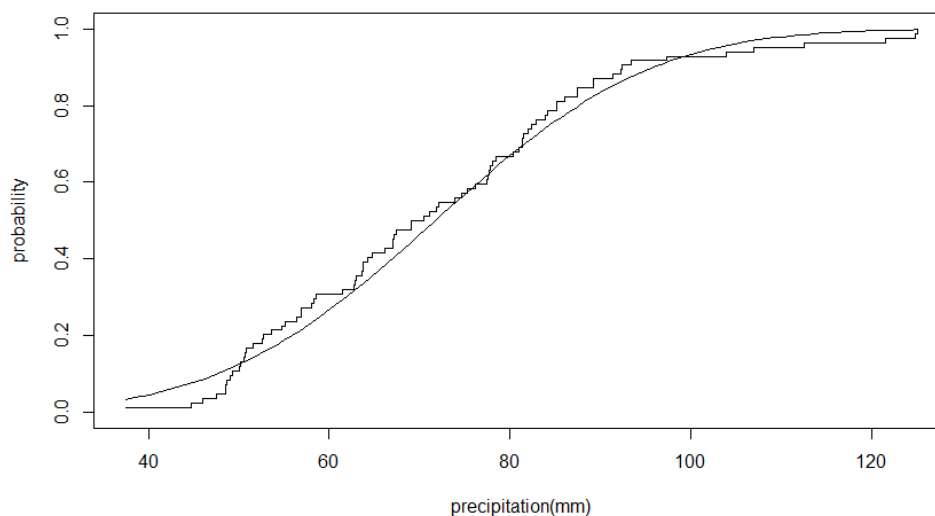
**2.c.** Tíðnirit með flatarmál sem er jafn einum. Teiknið ofan á það þéttifall normaldreifingar meðalgildi jafnt  $\bar{x}$  og staðalfrávik  $s$ .

**Svar:** Við notum skipunina `precipitationGraphs(2,FALSE)` og fáum:



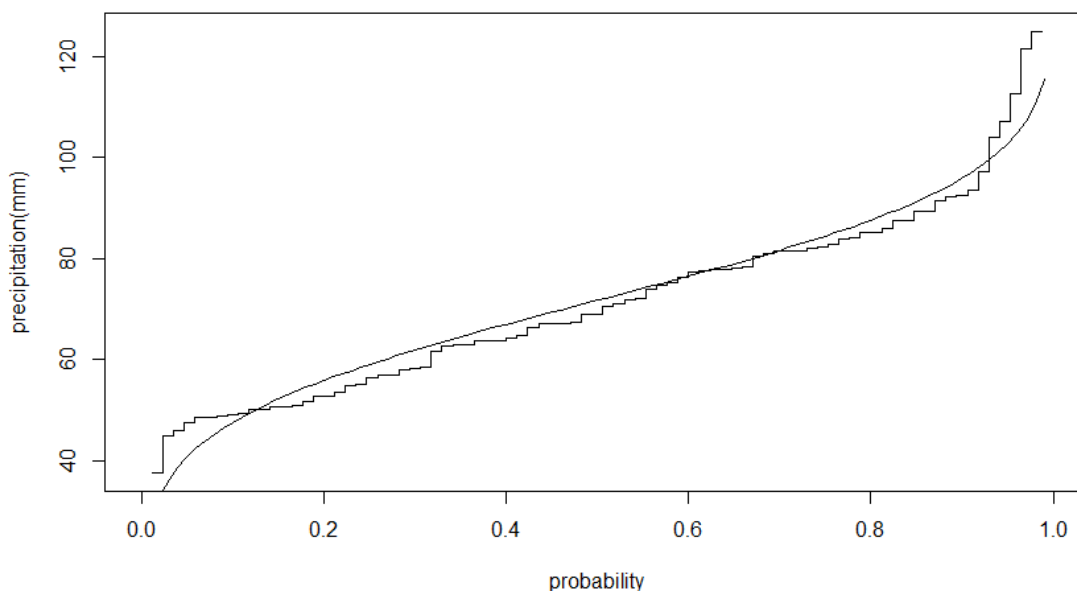
**2.d.** Dreififall úrtaks. Teiknið ofan á sætisfall úrtaksins sætisfall normaldreifingar með meðalgildi jafnt  $\bar{x}$  og staðalfrávik  $s$ .

**Svar:** Keyrum fallið aftur með skipun `precipitationGraphs(3,FALSE)` og fáum:



**2.e.** Sætisfall úrtaks. Teiknið ofan á það sätisfall normaldreifingar með meðalgildi jafnt  $\bar{x}$  og staðalfrávik  $s$ .

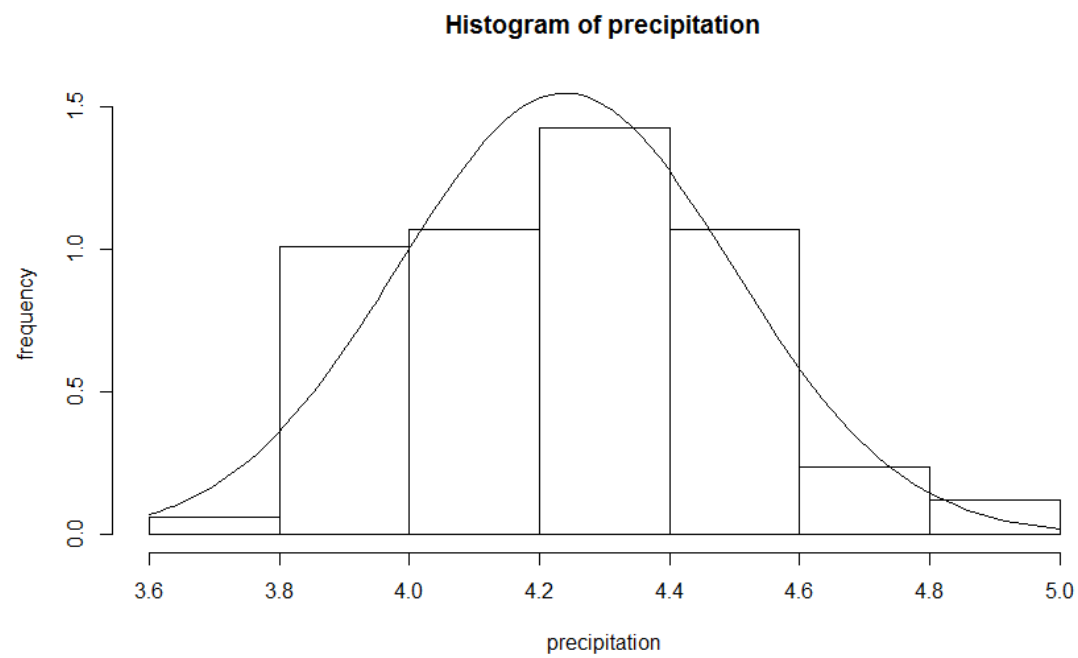
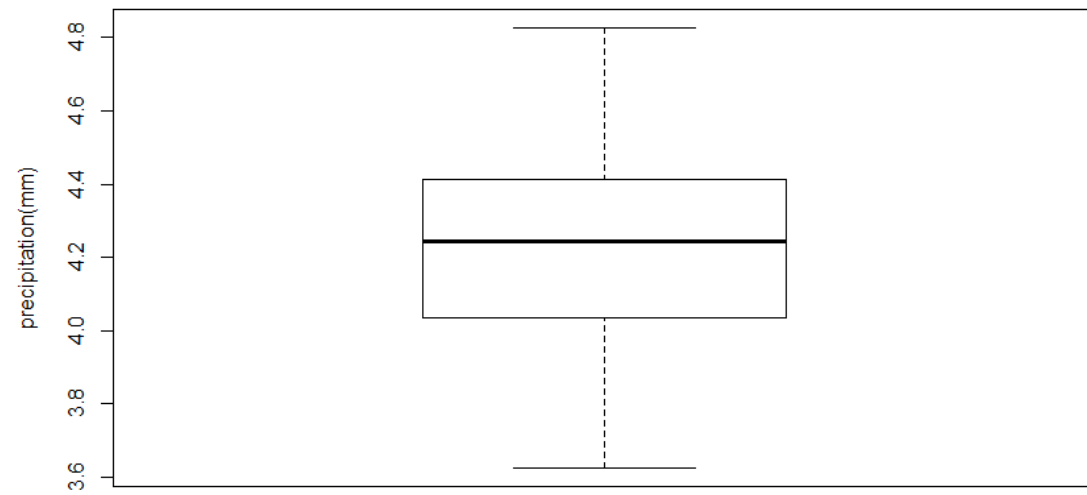
**Svar:** Keyrum þetta fall enn eina ferðina með skipun `precipitationGraphs(4,FALSE)` til að fá:

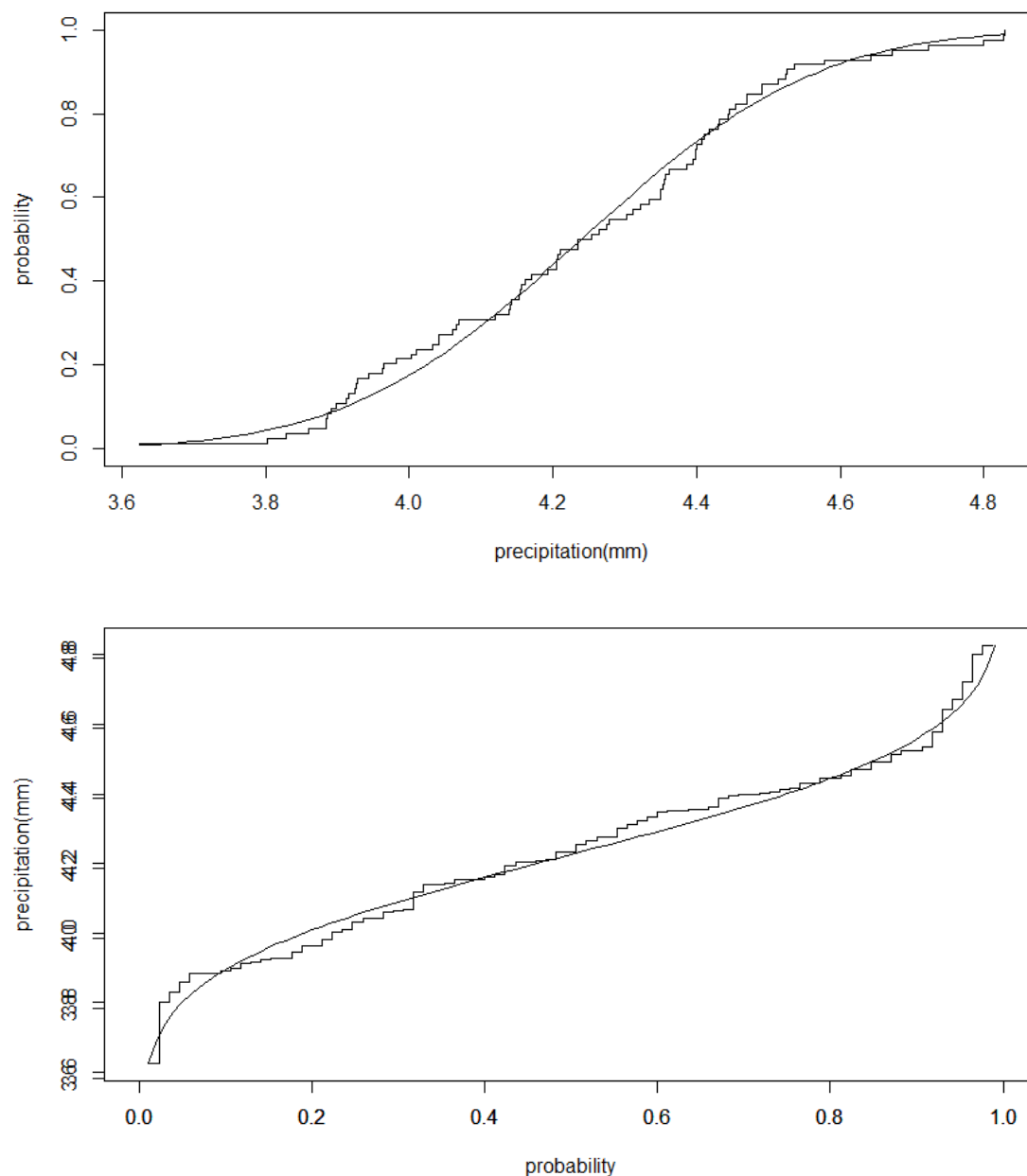


**2.f.** Teiknið sömu myndir og og í liðunum að ofan nema með logragildin af gögnunum.

**Svar:** Notum þetta mjög svo skemmtilega forrit nokkrum sinnum í viðbót nema hvað að núna er seinna viðfangið `TRUE` í öllum tilfellum. Þá fást þessar undurfögru og mjög svo fróðlegu gröf:







**2.g.** Lýsir normaldreifing gögnunum nægjanlega vel? Lýsir lognormaldreifing gögnunum nægjanlega vel? Notið myndirnar úr liðunum hér fyrir ofan til að rökstyðja svör ykkar.

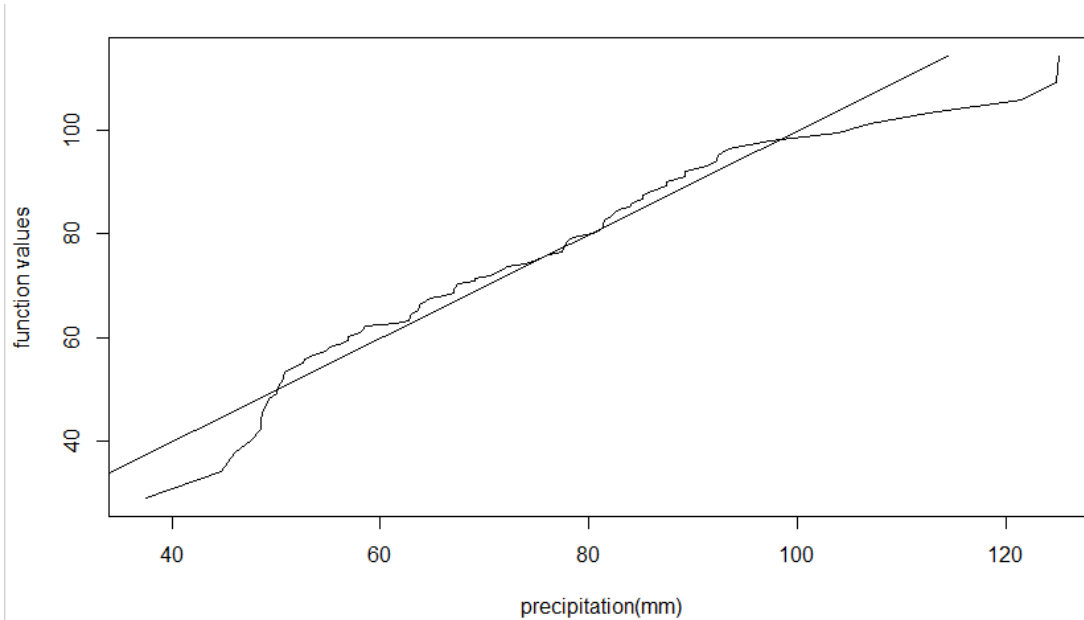
**Svar:** Bæði normaldreifing og lognormaldreifing lýsa gögnunum nokkuð vel eins og sést á ferlum að ofan (falla nokkuð vel að fræðilegum normaldreifingum). Hins vegar er lognormaldreifingin aðeins betri þar sem hún er nákvæmari og skekkjur eru minni.

### Verkefni 3

Hér á að teikna tvær myndir, þá fyrri byggða á gögnunum og hina byggða á logranum af gögnunum. Látum  $x_1 < x_2 < \dots < x_i$  tákna röðuðu gögnin

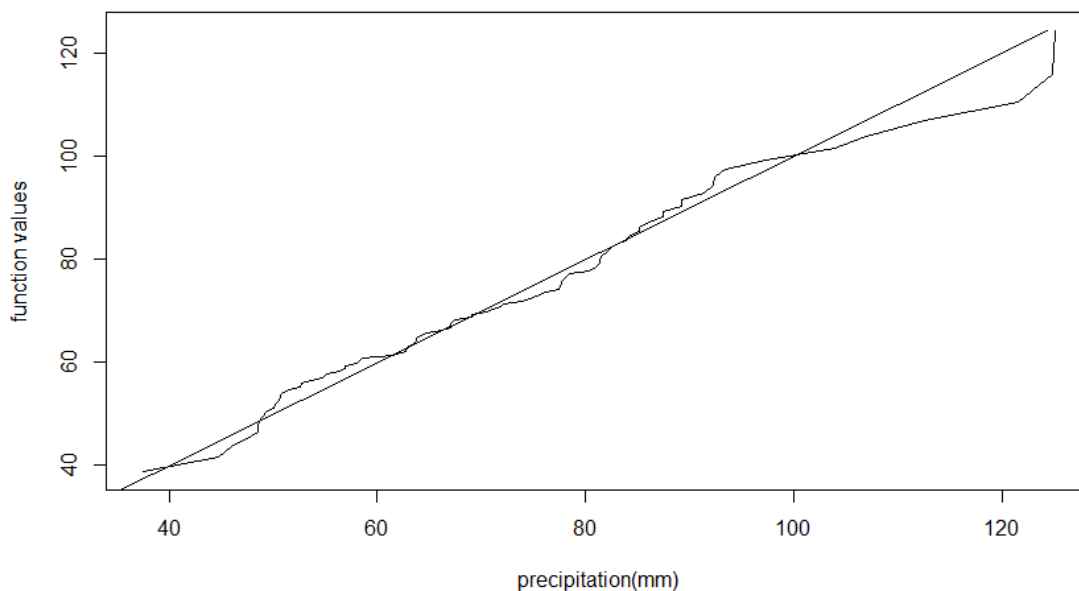
**3.a.** Teiknið  $x_i$  á x-ás og  $\mu + \sigma\Phi^{-1}(i/(n+1))$ . Teiknið á sömu mynd línuna  $y = x$ .

**Svar:** Keyrum fall með skipun `precipitationQnormGraphs(FALSE)` til að fá eftirfarandi graf:



**3.b.** Teiknið  $x_i$  á x-ás og  $\mu_y + \sigma_y\Phi^{-1}(i/(n+1))$ . Teiknið á sömu mynd línuna  $y = x$ .

**Svar:** Keyrum fall með skipun `precipitationQnormGraphs(TRUE)` til að fá eftirfarandi graf:



**3.c.** Því nær sem að punktarnir eru beinu línunni því betur lýsir dreifingin sem er lögð til gögnunum. Notið þessar tvær myndir til að svara eftirfarandi spurningum. Lýsir normaldreifing gögnunum nægjanlega vel? Lýsir lognormaldreifing gögnunum nægjanlega vel?

**Svar:** Af þessum tveimur myndum að dæma er log-normaldreifing þó nokkuð nákvæmari en normaldreifing. Hvort log-normaldreifing eða normaldreifing er "nægjanleg" fer alfarið eftir því hvað er verið að gera hverju sinni.



## Verkefni 4

Metið líkurnar á því að sólarhringsúrcoma á Fagurhólsmýri sé meiri en 115 mm, það er, metið  $P(X > 115)$

**4.a.** með því að finna hlutfall mælinga stærri en 115 mm.

**4.b.** með því að gera ráð fyrir að normaldreifingin lýsi mælingunum nægjanlega vel.

**4.c.** með því að gera ráð fyrir að normaldreifingin lýsi logranum af mælingunum nægjanlega vel.

**Svar:** Við keyrum seinasta fallið með kallinu `precipitationEstimation(115)` og fáum:

```
> precipitationEstimation(115)
Probability of more than 115 mm of precipitation built
on proportions: 0.0357142857142857
Probability of more than 115 mm of precipitation built
on normal distribution: 0.0108397295634248
Probability of more than 115 mm of precipitation built
on log-normal distribution: 0.024905613055092
```

**4.d.** Ef þið ættuð að meta  $P(X > 130)$ , hvaða aðferð munduð þið nota? Aðferðina í (a), (b), (c) eða einhverja aðra aðferð? Útskýrið val ykkar.

**Svar:** Af fyrrum liðum að dæma væri skynsamlegast að velja log-normal dreifingu. Hlutfalls-ágiskun er afar ónákvæm aðferð og byggist bara á gögnum sem eru þegar komin. Skekkjur í því módeli magnast upp með svona ágiskunum. Log-normaldreifing virðist draga mun frekar úr skekkjum heldur en normaldreifing og gefa betri niðurstöður. Eflaust eru til flóknari og betri aðferðir til að smíða svona líkön en áætla má að log-normaldreifing gefi nógu nákvæma mynd fyrir hin ýsmu verk.

# Viðauki - Forrit

## Verkefni 1

```

1 precipitationStats <- function() {
2
3   # Input of data for statistic processing
4   max_urkoma <- read.table("precipitation_fagurhm.txt")
5   urkoma <- max_urkoma[,2]
6   # Create matrix for storing statistics and label columns/rows
7   statMatrix <- matrix(rep(c(0),7),ncol=7,byrow=TRUE)
8   colnames(statMatrix) <- c("Mean", "Median", "Variance", "Standard
      Deviation",
9                               "First Quartile", "Third Quartile", "
      Interquartile Range")
10  rownames(statMatrix) <- c("Value")
11
12  # Calculate statistical values and input data
13  statMatrix[1,1] <- mean(urkoma)
14  statMatrix[1,2] <- median(urkoma)
15  statMatrix[1,3] <- var(urkoma)
16  statMatrix[1,4] <- sd(urkoma)
17  statMatrix[1,5:6] <- quantile(urkoma,probs = c(0.25,0.75))
18  statMatrix[1,7] <- statMatrix[1,6] - statMatrix[1,5]
19
20  # Export table to global scope for debugging
21  statMatrix <- statMatrix
22  print(statMatrix)
23
24  # Check if mean or median value is bigger
25  if(statMatrix[1,1] > statMatrix[1,2]){
26    message("Mean value is larger than median value.")
27  } else if(statMatrix[1,1] < statMatrix[1,2]){
28    message("Median value is larger than mean value.")
29  } else{
30    message("Mean value and median value are equal.")
31  }
32
33  # Calculate range for 50 percent of values with median values as
      middle point
34  message("Range with 50 percent of values with median as mid-point:
      [",
35          statMatrix[1,5], " ; ", statMatrix[1,6], "]"")
36
37
38  #####
39  # LOGARITHM STATISTICS
40  #####
41
42
43  # Calculate logarithmic values of data
44  logUrkoma <- log(urkoma)
45  # Create matrix for storing statistics and label columns/rows

```

```

46 logStatMatrix <- matrix(rep(c(0),7),ncol=7,byrow=TRUE)
47 colnames(logStatMatrix) <- c("Mean", "Median", "Variance", "
    Standard Deviation",
48                               "First Quartile", "Third Quartile", "
    Interquartile Range")
49 rownames(logStatMatrix) <- c("Logarithmic Values")
50
51 # Calculate statistical values and input logarithmic data
52 logStatMatrix[1,1] <- mean(logUrkoma)
53 logStatMatrix[1,2] <- median(logUrkoma)
54 logStatMatrix[1,3] <- var(logUrkoma)
55 logStatMatrix[1,4] <- sd(logUrkoma)
56 logStatMatrix[1,5:6] <- quantile(logUrkoma,probs = c(0.25,0.75))
57 logStatMatrix[1,7] <- logStatMatrix[1,6] - logStatMatrix[1,5]
58
59 # Export table to global scope for debugging
60 logStatMatrix <- logStatMatrix
61 print(logStatMatrix)
62 }

```

precipitationStats.r

## Verkefni 2

```

1
2 " Variable graphType is an integer and used for determining what
   method should be
3 implemented in creating the graph/plot (0 for scatter graph, 1 for
   box plot,
4 2 for histogram, 3 for cumulative distribution function or 4 for
   quantile function).
5 The variable logarithm is a boolean value to determine whether the
   logarithm of
6 the data should be used or not (TRUE means logarithmic values are
   used)."
7
8 precipitationGraphs <- function(graphType,logarithm){
9
10 # Input of data for statistic processing
11 max_urkoma <- read.table("precipitation_fagurhm.txt")
12
13 # Use logarithmic data if requested
14 if(logarithm){
15     max_urkoma[2] <- log(max_urkoma[2])
16     message("Notice: Logarithmic data being used")
17 }
18
19 # Often used statistics
20 sortedData <- sort(max_urkoma[,2])
21 meanOfData <- mean(sortedData)
22 stdDevOfData <- sd(sortedData)
23 n <- length(max_urkoma[,2])
24
25 # Draw scatter graph if required (graphType = 0)

```

```

26 if(graphType == 0){
27   plot(max_urkoma, xlab="year", ylab="precipitation(mm)")
28   # Check what year the precipitation was the most
29   maxRow = 1
30   for(g in 1:n){
31     if(max_urkoma[g,2] > max_urkoma[maxRow,2]){
32       maxRow <- g
33     }
34   }
35   message("Maximum precipitation was ", max_urkoma[maxRow,2], " in
        the year ", max_urkoma[maxRow,1])
36 }
37
38 # Draw boxplot if required (graphType = 1)
39 if(graphType == 1){
40   boxplotInstance <- boxplot(max_urkoma[2], ylab="precipitation(mm)
        ")
41   message("Number of outliers in boxplot: ", length(
        boxplotInstance$out))
42 }
43
44 # Draw histogram if required (graphType = 2)
45 if(graphType == 2){
46   precipitation <- max_urkoma[,2]
47   histogramInstance <- hist(precipitation, plot = F)
48   yLimit <- range(0, histogramInstance$density, dnorm(0)/sd(
        precipitation))
49   hist(precipitation, freq = F, ylim = yLimit, ylab = "frequency")
50   curve(dnorm(x, mean = mean(max_urkoma[,2]), sd = sd(max_urkoma
        [,2])), add = T)
51 }
52
53 # Draw cumulative distribution function if required (graphType =
        3)
54 if(graphType == 3){
55   # Cumulative distribution function for data
56   plot(sortedData, (1:n)/n, type = "s", ylim = c(0,1), xlab = "
        precipitation(mm)", ylab = "probability")
57   # Cumulative distribution function based on statistics from data
58   curve(pnorm(x, mean = meanOfData, sd = stdDevOfData), add = T)
59 }
60
61 # Draw quantile function if required (graphType = 4)
62 if(graphType == 4){
63   # Quantile function for data
64   sortedData <- sort(max_urkoma[,2])
65   p <- (1:n)/(n + 1)
66   plot(p, sortedData, type = "s", xlim = c(0,1), xlab = "
        probability", ylab = "precipitation(mm)")
67   # Quantile function built on data statistics
68   if(logarithm){
69     # Workaround for logarithmic graph
70     par(new = TRUE)
71     curve(qnorm(x, mean = meanOfData, sd = stdDevOfData), xlab = ""

```

```

    , ylab = "")
72 }else{
73   curve(qnorm(x,mean = meanOfData, sd = stdDevOfData), add = T)
74 }
75 }
76 }

```

precipitationGraphs.r

### Verkefni 3

```

1
2 # Parameter logarithm is a boolean variable to determine
3 # whether logarithmic data is used
4
5 precipitationQnormGraphs <- function(logarithm){
6
7   # Input of data for statistic processing
8   max_urkoma <- read.table("precipitation_fagurhm.txt")
9   log_urkoma <- log(sort(max_urkoma[,2]))
10
11   # Use logarithmic data if requested
12   if(logarithm){
13     message("Notice: Logarithmic data being used")
14   }
15
16   # Often used statistics
17   sortedData <- sort(max_urkoma[,2])
18   meanOfData <- mean(sortedData)
19   stdDevOfData <- sd(sortedData)
20   n <- length(sortedData)
21   logMeanOfData <- mean(log_urkoma)
22   logStdDevOfData <- sd(log_urkoma)
23
24   # Produce values according to assignment and plot them
25   if(logarithm){
26     functionValues <- exp(logMeanOfData + logStdDevOfData*qnorm((1:
27       n)/(n+1)))
28   }else{
29     functionValues <- meanOfData + stdDevOfData*qnorm((1:n)/(n+1))
30   }
31   plot(sortedData, functionValues, type = "n", xlab = "precipitation
32     (mm)", ylab = "function values")
33   lines(sortedData, functionValues)
34   lines(c(0,max(functionValues)), c(0,max(functionValues)))
35 }

```

precipitationQnormGraphs.r

### Verkefni 4

```

1
2 # Various estimations of input data

```

```
3
4 precipitationEstimation <- function(amount){
5
6   # Input of data for statistic processing
7   max_urkoma <- read.table("precipitation_fagurhm.txt")
8
9   # Often used statistics
10  sortedData <- sort(max_urkoma[,2])
11  meanOfData <- mean(sortedData)
12  stdDevOfData <- sd(sortedData)
13  n <- length(max_urkoma[,2])
14
15  # Find place in data where values are less than amount
16  place = 1
17  for(g in 1:n){
18    if(sortedData[g] > amount){
19      break
20    }
21    place <- g
22  }
23
24  # Calculate probability built on proportions
25  proportion <- (n - place)/n
26  message("Probability of more than ", amount,
27         " mm of precipitation built on proportions: ", proportion)
28
29  # Calculate the probability built on normal distribution
30  proportion <- 1 - pnorm(amount, mean = meanOfData, sd =
31                        stdDevOfData)
32  message("Probability of more than ", amount,
33         " mm of precipitation built on normal distribution: ",
34         proportion)
35
36  # Calculate the probability built on log-normal distribution
37  proportion <- 1 - pnorm(log(amount), mean = mean(log(sortedData)),
38                        sd = sd(log(sortedData)))
39  message("Probability of more than ", amount,
40         " mm of precipitation built on log-normal distribution: ",
41         proportion)
42 }
```

precipitationEstimation.r