

R Tutorial for NCCU-SAMSI Undergraduate Data-Science Workshop

Xinyi Li, SAMSI

03/04/2019

About R

R is a free software environment for statistical computing and graphics:

- a different implementation of S developed at Bell Lab;
- provides a wide variety of statistical and graphical techniques, and is highly extensible;
- open source;
- powerful IDE (integrated development environment), such as Rstudio.

Install R

1. Download the most recent version of R. The R FAQs and the R Installation and Administration Manual contain detailed instructions for installing R on various platforms (Linux, OS X, and Windows being the main ones).
2. Start the R program; on Windows and OS X, this will usually mean double-clicking on the R application, on UNIX-like systems, type “R” at a shell prompt.
3. As a first step with R, start the R help browser by typing `help.start()` in the R command window. For help on any function, e.g. the “mean” function, type `?mean`.

Install RStudio

1. Go to RStudio and click on the “Download RStudio” button.
2. Click on “Download RStudio Desktop.”
3. Click on the version recommended for your system, or the latest Windows version, and save the executable file. Run the .exe file and follow the installation instructions.

Data types

We can use variables without definition in advance.

Numbers

```
num = 3.14  
print(num)
```

```
## [1] 3.14
```

```
print(num + 1)
```

```
## [1] 4.14
```

```
print(typeof(num))
```

```
## [1] "double"
```

```
num.int = as.integer(num);
```

```
print(num.int)
```

```
## [1] 3
```

```
print(typeof(num.int))
```

```
## [1] "integer"
```

We can use R as a calculator, e.g. $2 * 2$, $\log(2)$, $\sqrt{2}$, 2^3 .

```
x = 2
```

```
print(x * 2)
```

```
## [1] 4
```

```
print(log(x))
```

```
## [1] 0.6931472
```

```
print(sqrt(x))
```

```
## [1] 1.414214
```

```
print(x ^ 3)
```

```
## [1] 8
```

```
print(x ** 3)
```

```
## [1] 8
```

Data frame

```
y = 10:12
```

```
print(y)
```

```
## [1] 10 11 12
```

```
z = c(1, 3, 5)
```

```
print(z)
```

```
## [1] 1 3 5
```

```
print(z[1])
```

```
## [1] 1
```

```
df = data.frame(y = y, z = z)
```

```
print(df)
```

```
##      y z
```

```
## 1 10 1
```

```
## 2 11 3
## 3 12 5

print(class(df))

## [1] "data.frame"

print(df$y)

## [1] 10 11 12

print(df$z)

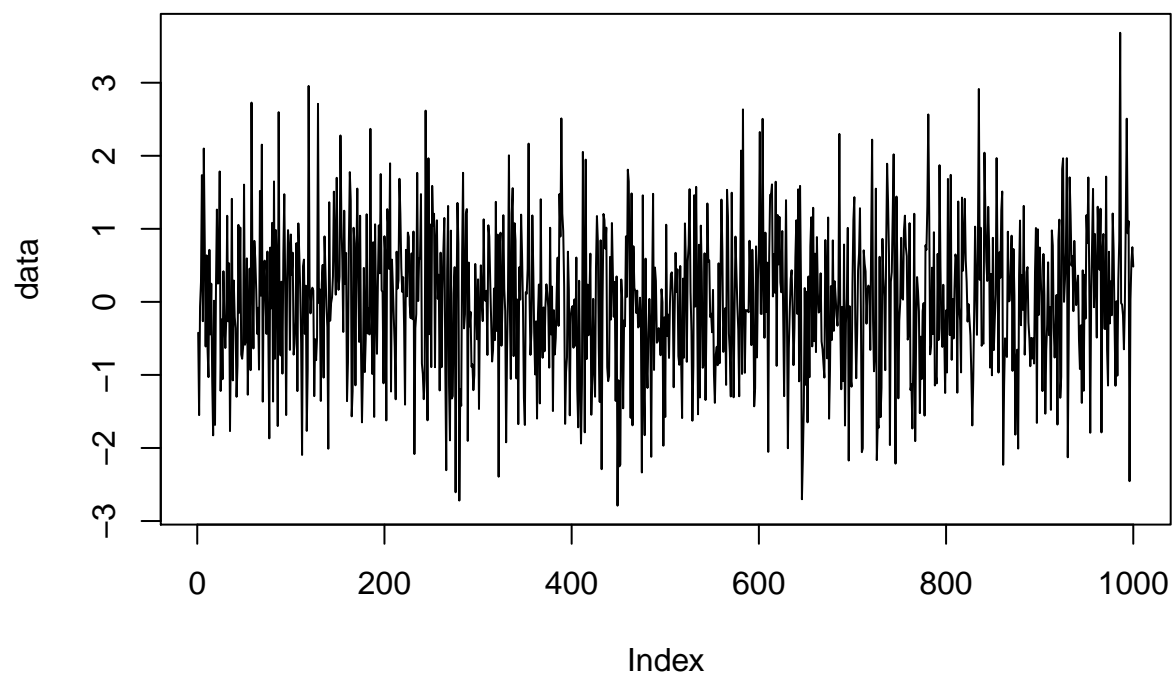
## [1] 1 3 5
```

Exercise: Create a data frame containing name, gender, grades, etc.

Basic plots

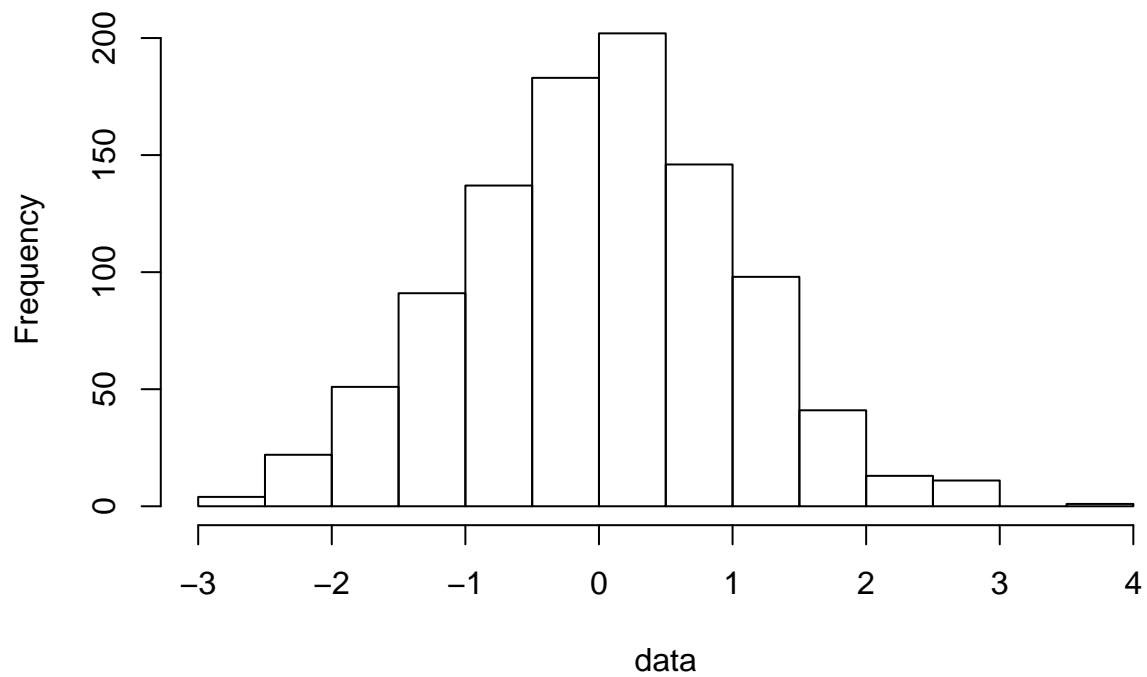
Use of “hist” function

```
set.seed(2018)
data = rnorm(1000)
plot(data, type = 'l')
```



```
hist(data)
```

Histogram of data



Use of “plot” function

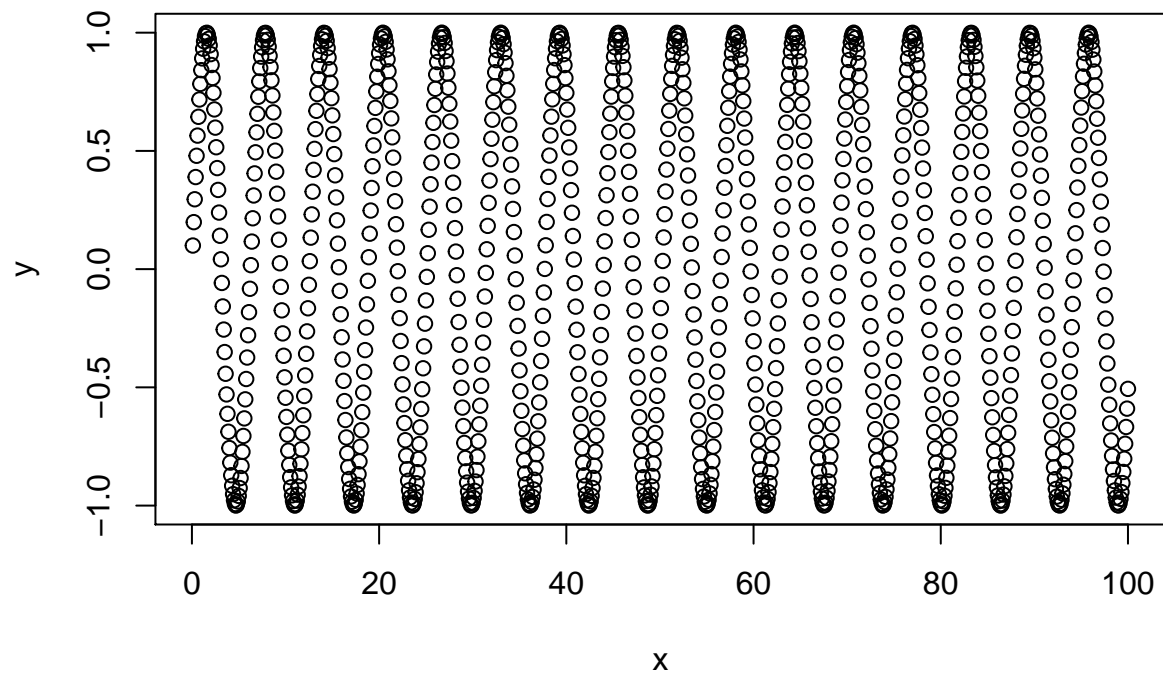
```
x = 1:1000/10  
y = sin(x)  
print(head(x))
```

```
## [1] 0.1 0.2 0.3 0.4 0.5 0.6
```

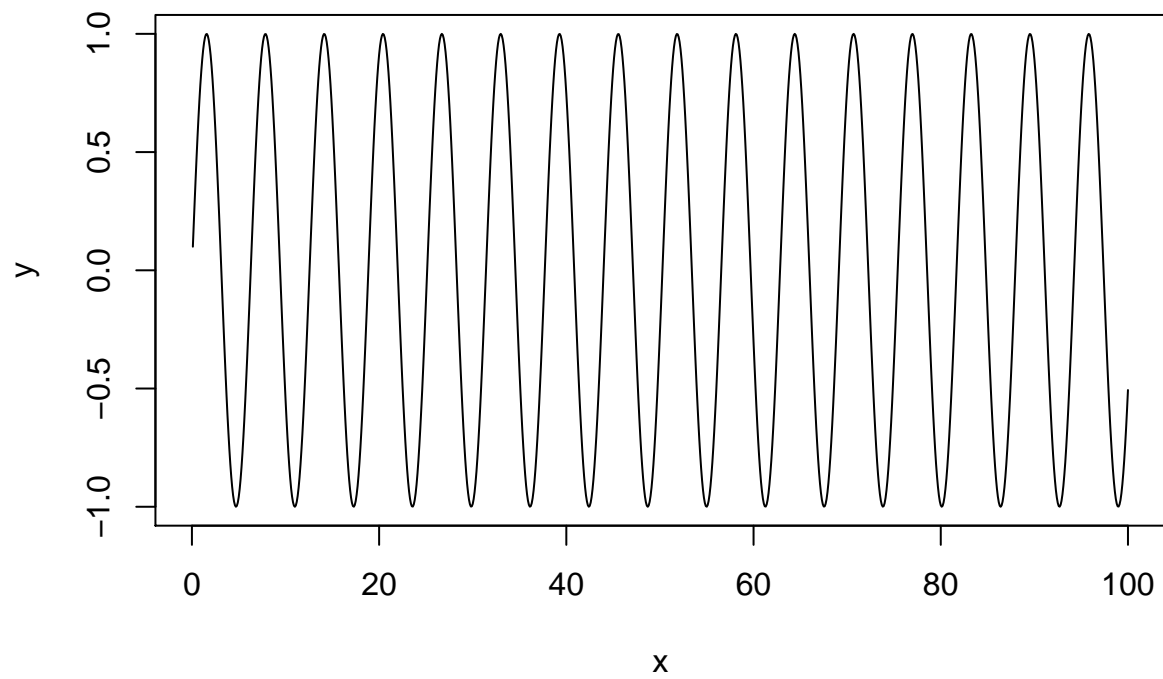
```
print(tail(y))
```

```
## [1] -0.8577953 -0.8021964 -0.7385822 -0.6675884 -0.5899242 -0.5063656
```

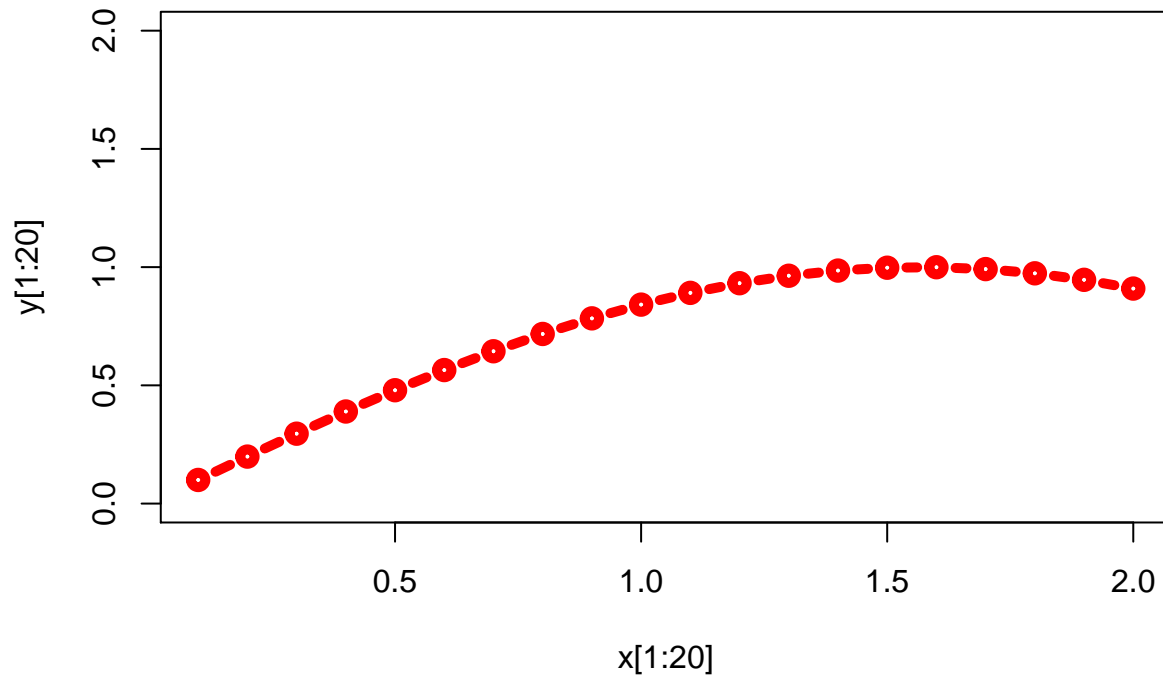
```
plot(x, y)
```



```
plot(x, y, type = "l")
```

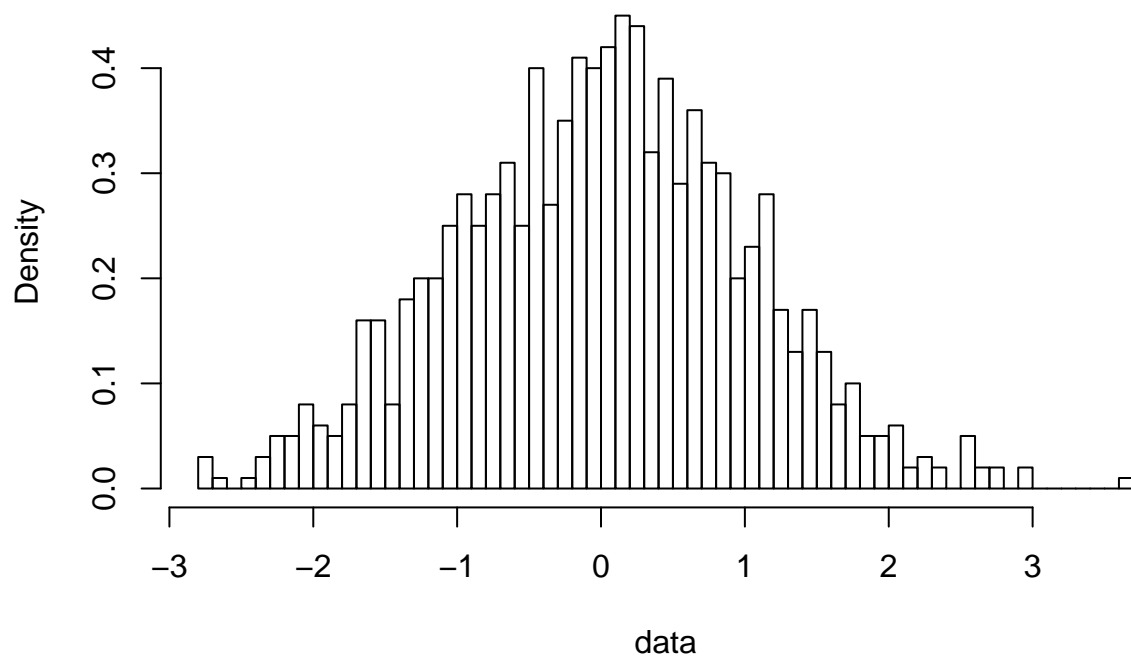


```
plot(x[1:20], y[1:20], type = 'b', col = 'red', lwd = 5, ylim = c(0, 2))
```



Exercise: Use `hist()` and change parameters to generate the figure as follow.

Histogram of data



Importing Data

Download data from https://github.com/LiXinyi/SAMSI_Diversity_Workshop/blob/master/CanadianWeather_month.csv. Original data are available at R package `fda`.

- Importing a single file

```
dat = read.csv("CanadianWeather_month.csv", header = TRUE)
print(class(dat))
```

```
## [1] "data.frame"
```

```
print(dim(dat))
```

```
## [1] 420 4
```

```
print(head(dat))
```

```
##      Temp  Precip Month  Region
## 1 -4.654839 4.651613   Jan St. Johns
## 2 -5.325000 4.735714   Feb St. Johns
## 3 -2.532258 4.235484   Mar St. Johns
## 4  1.256667 3.616667   Apr St. Johns
## 5  5.793548 3.251613   May St. Johns
## 6 10.786667 3.270000   Jun St. Johns
```

```
print(names(dat))
```

```
## [1] "Temp" "Precip" "Month" "Region"
```

```
print(table(dat$Month))
```

```
##
## Apr Aug Dec Feb Jan Jul Jun Mar May Nov Oct Sep
## 35 35 35 35 35 35 35 35 35 35 35 35
```

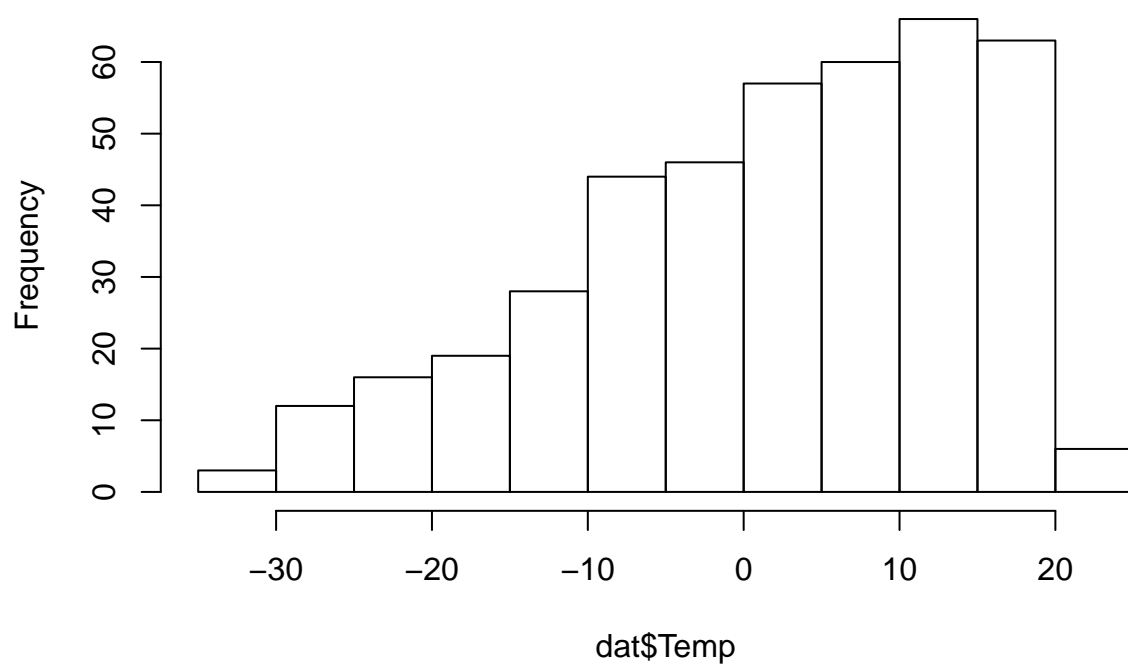
```
print(table(dat$Region))
```

```
##
##      Arvida Bagottville      Calgary Charlottvl  Churchill      Dawson
##      12      12      12      12      12      12
##      Edmonton Fredericton      Halifax      Inuvik      Iqaluit      Kamloops
##      12      12      12      12      12      12
##      London      Montreal      Ottawa Pr. Albert Pr. George Pr. Rupert
##      12      12      12      12      12      12
##      Quebec      Regina      Resolute Scheffervll  Sherbrooke      St. Johns
##      12      12      12      12      12      12
##      Sydney      The Pas  Thunderbay      Toronto Uranium Cty  Vancouver
##      12      12      12      12      12      12
##      Victoria Whitehorse      Winnipeg      Yarmouth Yellowknife
##      12      12      12      12      12
```

Exercise: Use `hist()` and `plot()` to get basic idea of the data.

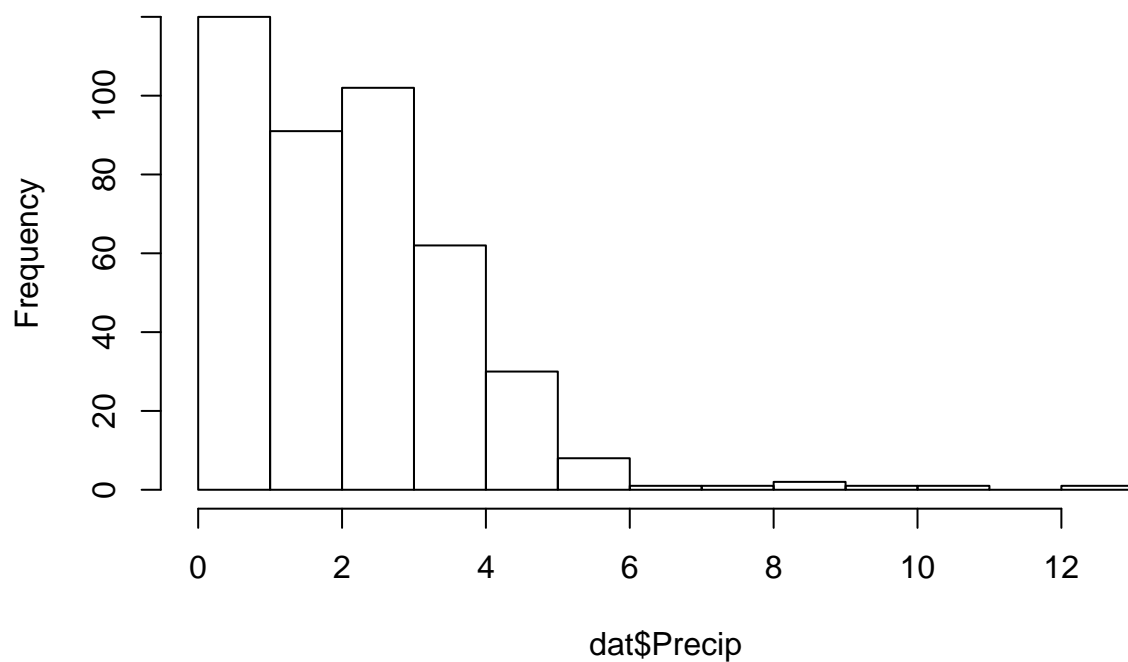
```
hist(dat$Temp)
```

Histogram of dat\$Temp

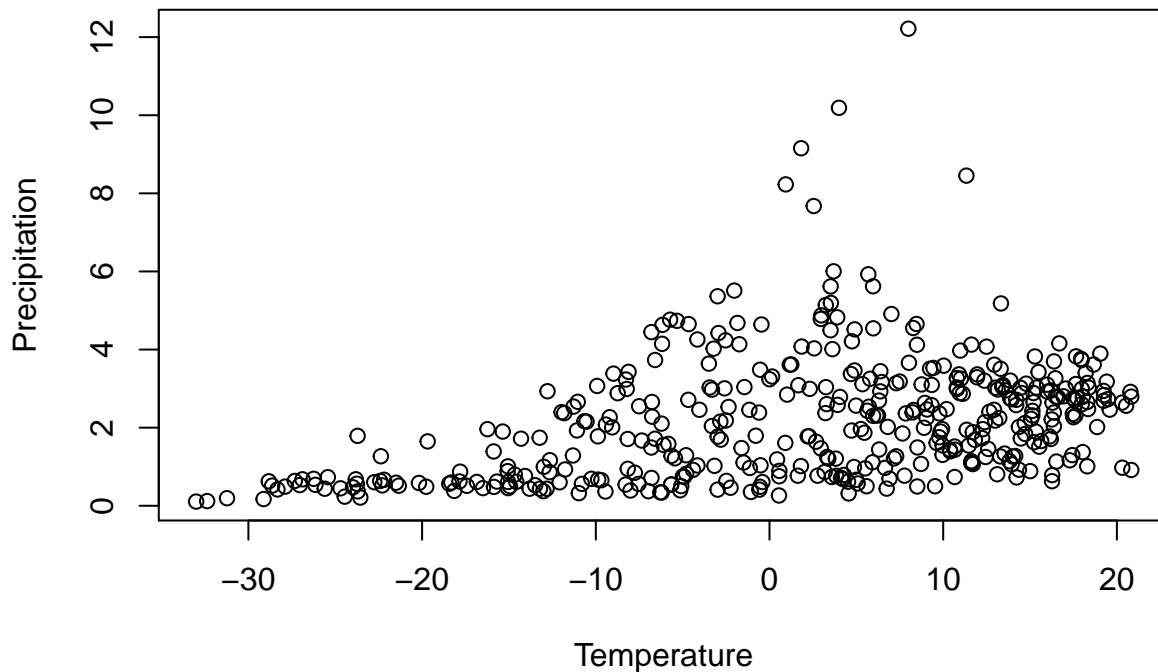


```
hist(dat$Precip)
```

Histogram of dat\$Precip



```
plot(dat$Temp, dat$Precip, xlab = "Temperature", ylab = "Precipitation")
```

Packages

How to install a package

Install from source

Download the add-on R package, for example, “fda”, put it in the directory “/data/Rpackages”, and install the package using the command:

```
install.packages("fda", lib = "/data/Rpackages")
```

Install from repository

Vast array of packages are available at the Comprehensive R Archive Network (CRAN) and BioConductor repositories. Both CRAN and BioConductor are open source, well structured, tested and operating. While both repositories provide abundant packages covering various data analysis tasks, BioConductor is more focused on providing tools for the analysis of high-throughput genomic data. In addition, there are slight differences in the command for package installation.

- Install from CRAN (e.g. R package “fda”):

```
install.packages("fda", repos = "http://cran.us.r-project.org")
```

- Install from BioConductor (e.g. R package “dada2”):

```
## try http:// if https:// URLs are not supported
source("https://bioconductor.org/biocLite.R")
```

```
## Bioconductor version 3.7 (BiocInstaller 1.30.0), ?biocLite for help
```

```
## A newer version of Bioconductor is available for this version of R,
```

```
## ?BiocUpgrade for help
```

```
biocLite()

## BioC_mirror: https://bioconductor.org
## Using Bioconductor 3.7 (BiocInstaller 1.30.0), R 3.5.1 (2018-07-02).
## Old packages: 'backports', 'BH', 'bookdown', 'broom', 'callr', 'class',
## 'clipr', 'codetools', 'colorspace', 'curl', 'data.table', 'dbplyr',
## 'devtools', 'digest', 'dplyr', 'DynTxRegime', 'evaluate', 'fansi',
## 'forcats', 'ggplot2', 'git2r', 'grpreg', 'haven', 'httpuv', 'httr',
## 'igraph', 'jsonlite', 'knitr', 'later', 'lattice', 'markdown', 'MASS',
## 'Matrix', 'mgcv', 'mime', 'modelObj', 'modelr', 'openssl', 'pillar',
## 'pracma', 'processx', 'ps', 'purrr', 'R6', 'RandomFields',
## 'RandomFieldsUtils', 'Rcpp', 'RcppEigen', 'RcppParallel', 'readr',
## 'readxl', 'rlang', 'rmarkdown', 'rstudioapi', 'spam', 'stringi',
## 'stringr', 'survival', 'tibble', 'tidyr', 'tinytex', 'tseriesChaos',
## 'vegan', 'xfun'

biocLite("dada2")
```

How to load functions from a package

Type the following command in R console to load the package.

```
library(fda)

## Loading required package: splines
## Loading required package: Matrix
##
## Attaching package: 'fda'
## The following object is masked from 'package:graphics':
##
##      matplot

library(dada2)

## Loading required package: Rcpp
attach(CanadianWeather)
names(CanadianWeather)

## [1] "dailyAv"      "place"        "province"     "coordinates"
## [5] "region"       "monthlyTemp"  "monthlyPrecip" "geogindex"
```

Exercise: Explore by yourself for the Canadian Weather data.