

Nutrition & Planetary Health Research Group

Introduction to R

Zentrum für Entwicklungsforschung (ZEF)

Author: Mahir Bhatt

2025/2026

Difference between R & RStudio



(a) R is a programming language



(b) RStudio is programming environment

Working with RStudio

- Always create a project for specific analysis.
- Set the project directory to the location where your data is stored.
- Do not open R files from multiple directories at the same time (**RStudio could crash!**)

Using Rmarkdown vs RScript

There is nothing wrong with using either of them, however .Rmd has bit of an advantage.

- .Rmd helps you make your code more dynamic by combining text with code and output.
- You can knit .Rmd files in to html, pdf or word format
- .Rmd ensures that your analysis is reproducible.
- You can include interactive elements like shiny apps & widgets with .Rmd
- You can easily make your code accessible on Github with .Rmd

Packages in R

- There are pre-installed packages in R which can perform basic statistical functions.
- To perform specific tasks, you would have to install package the contains the given function
- R has a large repository of packages (**Constantly updated!**)

Packages in R

- There are pre-installed packages in R which can perform basic statistical functions.
- To perform specific tasks, you would have to install package the contains the given function
- R has a large repository of packages (**Constantly updated!**)
- Example of a mathematical operation

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Lists, Vectors & Data Frames

Lists

It is the most basic type of data storage unit in R. It generates a list of inputs stored which can be heterogeneous (It is like storing different attributes of a person in a single variable).

Vectors

Vector is a simple one-dimensional array of data. Elements of all vectors are homogeneous (all numeric, all characters). It is more like storing the same attribute of a different person.

Data Frames

They are the building block of data analysis in R. It stores data in tabular format (two-dimensions). Each column of the table acts as an individual vector.

Loading a data frame in R

- R is capable of loading any type of data frame (SAS, STATA, SPSS, csv, excel, *QGIS*, *ArcGIS*)
- Unless you are loading a dataset from another which has an extension of another statistical software, always use **csv**!

Visualizing data in R

R has inbuilt capabilities to graphically present your data. Below are examples of some packages that can be used.

- `ggplot2` → Combine multiple data sources to create graphical summary
- `plotly` → Create interactive plots in 2D and 3D
- `Leaflet` → Build interactive maps
- `tmaps` → Use shapefiles to produce descriptive maps
- `highcharter` → Used for dynamic charting
- `Lattice` → Powerful tool for visualizing multivariate data

We will first look at the use of `ggplot2` and slowly advance skills with other packages.

Types of ggplots

```
1 # geom_boxplot()
2 ggplot(df1, aes(x=, y
   =)) +
3   geom_boxplot() +
4   theme_bw() +
5   labs(x = "", y = "",
        title = "")
6
7 # geom_smooth()
8 ggplot(df1, aes(x=, y
   =)) + geom_smooth() +
9   theme_bw() + labs(x = "",
   , y = "", title = "")
```

```
1 # geom_line()
2 ggplot(df1, aes(x=, y
   =)) +
3   geom_line() +
4   theme_bw() +
5   labs(x = "", y = "",
        title = "")
6
7 # geom_point()
8 ggplot(df1, aes(x=, y
   =)) +
9   geom_boxplot() +
10  theme_bw() +
11  labs(x = "", y = "",
        title = "")
```

Types of ggplots

- ggplot allows you to create wide range of plots to graphically show your data, you do make more plots your skills with automatically improve.

Some examples

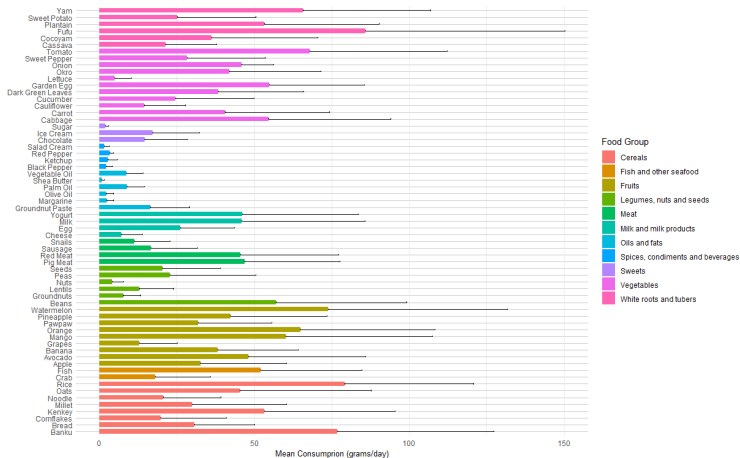
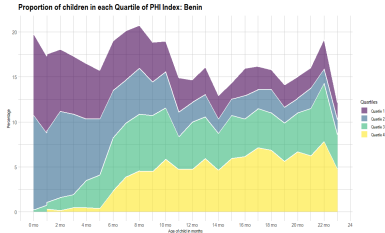
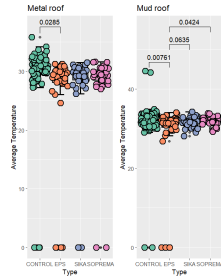


Figure: Mean and standard deviation of all the food items consumed by children in Agogo

More examples



(a) Proportion of children in each quartile of Planetary Health Diet Index by age



(b) Tukey plot of show effect of cool roof paints on indoor temperature in Mud and Metal roof

Linear models in R

Beginning with bivariate data exploration

```
1      #make a plot looking
2      at two variables
3      plot(df1$height, df1$
4      weight)
5
6      #Cross tabulation can be
7      used for categorical
8      data
9      table(df1$rdt, df1$sex,
10      useNA="always")
```

```
1      # variance, correlation
2      and covariance
3      var(df1$height
4      cor(df1$age, df1$
5      height)
6      var(df1$height,
7      df1$weight)
```

Fitting a linear model

- Let us consider the simple case of linear relationship between a response variable and predictor variable
 $y \approx \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \delta$ Where
- β_0 is intercept
- β_1, \dots, β_n are coefficients
- δ is error term

```
1  #simple linear model in R
2  model1<-lm(hemoglobin ~ weight, data=df1)
3  summary(model1)
4
5  model2<-lm(hemoglobin~weight+rdt, data=df1)
6  summary(model2)
```

Tables

I think this is enough for today. Thank you for
your attention 😊