

→ imbalance vectors :- 2 vectors of imbalance length can be added.  
 → the shorter one repeats itself and adds it to the larger one until the larger one is completed.

→  $[x \leftarrow val]$  - vector with one element

- $c(values) = \underline{\text{vector}}$  (which is 1-D) - grp of elements
- only one type of data can be their in a vector.  
 ↳ if numeric and character type are in vector then vector will convert the numeric one to character and the data type of whole vector would be character.

- In R, indexing starts from 1  
 ↳ comma lagana jaroori. → Container [nul-change] [after = 2] - 2<sup>nd</sup>/index
- List → collection of data of different DT - Just like Python - It is a 2-D Data Structure.  
 ↳  $list(values)$ .      \*  $\star [c(l_1, l_2) - \underline{\text{Merge 2 list}}]$

→ we can provide alias to the values of a list  $\star [append(l_1, l_2)]$

$var \leftarrow list(name = "Gaurav", age = 18) | var['name'] = after$   
 ↳ alias → Stores element as container  $= index$

★ Container → Stores element as container | can also use  $names()$   
 # To remove element from list → Assign it as NULL  
 $list[2] \leftarrow NULL$        $\star [\$ - \underline{\text{Access}}]$

- We can have list and it can have a vector as one of its member → NESTING.

Providing alias in vector →  $names()$

•  $list[index] \leftarrow$  new value

$vec \leftarrow c(1, 2, 3, 4)$

→ Providing alias to elements of vector.

↳ manipulating list

•  $names(vec) \leftarrow c('one', 'two', 'three', 'four')$  Elements.

•  $vec['two']$  is equal as calling  $vec[2]$ .

unlist() → Converting list to a vector

[concept of Recycling] → Recycling

sort(vec) - sort(vec, decreasing = TRUE).

length(vec)

— x — x — x — x — x —

• Lowest/min and Highest/max value of a vector

→ range(vec)

— x — ( ) — x — — — x —

• Generating sequence of numbers.

→ seq(1, 100, by = 10)

↳ skipping steps

• mean - sd - var (variance)  
↳ std deviation ( $\sqrt{var}$ )

$in = \text{Python}$   
 $.in[1] = R$

mean(vec)

sd(vec)

var(vec)

"ele" · 1 · in · 1 · list.

$.in[1] = R$

• To remove an element from vector.

INDEX

→ Just call negative of the index we want to ~~remove~~ Remove.

→ if I have to remove vec[1] element,

then Just call vec[-1] [It will remove that element vec[1]].

— x — x — x — x —

• Just call negative of the index we want to Remove.

Matrix → 2D [ vector, nrow = , ncol = ] parameters.  
Data for matrix

`mat → matrix(c(1,2,3,4), nrow=3, ncol=3, byrow=TRUE).`

Array → Multi-Dimensional Data structure.

$\rightarrow [\text{Data}, \dim]$  [ $\dim$  as vector] [ $\dim \leftarrow c(a, b, c)$ ]  
30 or more.

array → array ( $c(1, 2, 3, 4)$ ,  $\dim \leftarrow c(3, 3, 4)$ )

- Indexing is same as in python  $\rightarrow [os[r, c]] / [os[r_1:r_2, c_1:c_2]]$

Factors → Categorize the Data, and Store Pts as levels.

- They can store Both Strings and integers
  - Useful In columns which have a limited number of unique values.

→ To convert use `factor(Data structure)`

- Categorize + Store as levels → identifies unique values and store it as levels.

Example → `data <- c("East", "West", "East", "North", "North", "East", "West", "West", "East", "North").`

```
data-fac <- factor(data)
```

cat - Stores unique values as levels.

- `data-jac` → East, North, West
- `Summary(data-jac)` → East, North, West  
 4      3      4 → categorize the Data\*

For Dataframes [ To create `Dataframe` → `data.frame()` ]

- 1) Create features using data structures ↪ (vec, list etc.)
- 2) Create Dataframe → `data.frame(feat1, feat2)`

↑ o

`data.frame(feat1, feat2, stringAsFactors = TRUE)`

### Imp

- Generally / usually / in prov version when we create a col / feature of text or categorical col in dataframe → factors of it are automatically created BUT

⇒ Now we have to specify it by using `stringAsFactors`.

Refering Columns → `data.frame['col name']` / `data.frame$colname`

F - strings in R → library(glue)

`glue("Addition of {a} and {b} is {a+b}")`

## Functions in R.

Keyword

→ `function-name <- function ( )`  
- parameters

→ `add <- function ()`  
No parameters

eg. `Add <- function ()`

↳

`a = 5`

`b = 10`

`print(glue("Addition of {} and {} is {}"))`

↳

————— x ————— x ————— x ————— x —————

★ We can create global variables INSIDE A FUNCTION in R.

`var <- value`

`add-param <- function(a, b=10)`

↳ `global-var <- 10` → global var declaration inside a funcn.  
`return (a+b)`

↳

————— x ————— x ————— x —————

Type Conversion → as.desired.output.Data-Type()

`c <- 9.5` ] numeric - default data type for numbers.  
`d <- 60`

`c <- as.integer(c)`

`d <- as.integer(d)`

To check  
data type →  
`class()`

Imp ]

- Integers can be converted into characters ( $65 \rightarrow "65"$ ).
- Reverse is not True (if character is a alphabet).

Characters can be converted to integers only if that character is a number.

(Just like python).

— X — X — X —

Try-Catch [exception / Error Handling].

→ Syntax → try catch ( )

try catch(operation, error = function(e)  
print()).  
For errors

try catch(operation, warning = function(e)  
print()).  
for warnings

ch = 'a' | ch + var - error  
var = 5 | as.integer(ch) - warning | .

— X — X — X —

READING FILES IN R.

1) CSV files → read.csv()

2) Excel files →  
1) install package  
2) load that package (library())  
3) use it

- install.packages('readxl')
- library(readxl)
- Syntax → read-excel().

3] Text files → `readlines()` [var ← readlines()]  
 ↳ reads line by line

`len(var)` → Number of lines in txt files  
`nchar(var)` → Number of characters in each line.

# Reading char by char in txt file → scan()

`var1 ← scan("pathoffile", "")`.

— x — x — x —  
**WRITING and Saving into files.**

## 1] Txt

`Data ← matrix(c(1,2,3,4), nrow = 3, ncol = 4, byrow = TRUE)`.

write(Data, file = 'filepath', ncol = 5, sep = ',')  
 ↳ Irrespective of Data - In how many columns data should be saved in txt file.  
 ↳ Sep B/w values

2) into csv

dataframe

`write.csv(df, file = "path", row.names = FALSE)`

↳ as in `df` we already have col names.

3) in excel

- 1) install package - `xlsx`
- 2) load - `library(xlsx)`
- 3) use - `write.xlsx`

→ `install.packages('xlsx')`  
→ `library(xlsx)`  
→ `write.xlsx(df, file = 'path', Sheetname = 'Sheet1', col.names = TRUE, row.names = FALSE)`.

---

X — X — X —  
Basic String Operations in R.

# 1 → Number of characters in a string (= or -)

→ `nchar(var)`

# 2 → Convert to 'lower' case and 'upper case'

→ `tolower(var)`  
`toupper(var)`

## # 3 To replace characters in a String

Chartr (' ', '.', 'string')

[r → to replace]

• what to replace

space

• with what

Dot

every space in  
string will be  
replaced by Dot.

## # 4 String Split

→ 1<sup>st</sup> → strsplit

l To get individual  
words.

• 1<sup>st</sup> → strsplit

['Death and freedom', '']  
[by / from where]  
[wherever we have a space - split the  
content]

• 2<sup>nd</sup> → and then unlist those words using unlist

eg → catch ← strsplit("Death and freedom", '')  
words ← unlist(catch)

## # 5 Sub-String (start = , end = )

→ Substr("Death and freedom", start = 1, end = 7)

[Indexing in R starts  
from 1]

#6 - trim\_ws → Removes white space from  
white space front and end of the string.

X X X  
Dates in R YYYY-MM-DD  
(Year-Month-Date)

1) UNIX (from 1970-01-01) To date format

var ← as.POSIXct (col, origin = '1970-01-01')  
result ← as.Date(var)  
Date class object [YYYY-MM-DD]

etc ← as.Date (col, '·1·Y/·1·M/·1·d')

last / This argument tells  
to expect the input in this slash  
format. (and Y, m, d).

→ ·1·Y → year (4-Digit) Representation

→ ·1·y → 2-Digit year Representation

→ ·1·M → Month Representation

→ ·1·B → FULL Month name

→ ·1·b → Three character abbr month name

• 1. A = full weekday Name.

• 1. a = 3 char abbr weekday Name.

String to date  $\rightarrow$  as.Date()

as.Date("1994/06/07", "%Y/%m/%d")

→ We can use comparison and arithmetic operators B/W Dates. (subtraction /  $>$ ,  $<$ ,  $=$ ,  $\leq$ ).

# useful date functions.

→ sys.Date()  $\rightarrow$  current date as R date object.

→ date()  $\rightarrow$  current DATE and TIME

# getting what we want from DATE

weekdays(sys.Date())

months(sys.Date())

quarters(sys.Date())

# Sequence of Dates

seq (sys.Date(), by = 'day', length.out = 10)

10 Dates, with ↑ by single Day

by How much like by day increase pt.

### Regular expressions in R

- match pattern and strings in text.

# . [period] = matches with any character

# [+] → matches the preceding Pattern Element  
One or more time

# [\ ] → ensures period matches w/ period  
and not any character.

# grep ('@+[ ]+', c(values), value = TRUE) → gives values that matches RE

→ To get result (extract the RE just not values).

result ← regexpr ("@+[ ]+", c(values))

regmatches (c(values), result).

↳ gives out the extracted RE

o/o pipe/o

## More Matrices

- Adding more rows and cols →  $\text{rbind} / \text{cbind}$  ↗ matrix combine

Removing Rows and cols

$\text{rbind}(\text{matrix}, \text{c}(\text{values}))$

~~$\text{matrix}[-\text{RowIndex}, -\text{ColIndex}]$  (removes whole Row and col)~~

~~$X$~~

~~$X$~~

~~$X$~~

Dplyr

$\%>%$  → pipeline →  $\text{function}(y) = \text{function}(n, y)$

Select() → mutate() → Column FUNC

filter() → distinct() → arrange() → Row FUNC

group\_by() → summarise() → Count() → Group func

goes Together

Select() → Select Columns by Name or Helper func.

$df \cdot | \cdot > \cdot \cdot \cdot \text{Select}(\text{col1}, \text{col2}, \text{col3}, \text{col4}, \text{new col name} = \text{col5})$

$df \cdot | \cdot > \cdot \cdot \cdot \text{Select}(\text{col1} : \text{colm})$

includes n and m

$df \cdot | \cdot > \cdot \cdot \cdot \text{Select}(-\text{col1}, -\text{col2})$

will not show col1 and 2 but every other

→ one-of  $(c(" ", " "))$  → starts-with (" ")

• Select()

→ Contains("") → Select columns whose name contains Spec String

→ ends\_with("")

every thing() ←  
every column

matches ("!t.") ←  
name matches with  
REGULAR EXP

mutate() → Compute or Append one or more new columns.

- `grep(" ")` → matches pattern and test within  
  ↳ gives True and/or False

mutate(new\_col\_name = cond^n on which col name is decided)

↳  
`mutate(is_collab = grep("featuring", artist))`   %>%  
                  ↳ col name select()

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

filter() → Extract Rows that meet logical criteria

`select() %>% filter(col-name and cond^n)`   [ ] for mul col

- arrange → sorting or viewing data order wise

↳  
`arrange(data) | arrange(desc(data))`

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

- distinct() → To remove duplicate data (rows)  
  ↳ distinct(rows, name)

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

Group\_by (group\_by) → summarise(max)

max,  
mean,  
var,  
std

( )

\_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_ X \_\_\_\_\_

Count() → Count number of rows with unique  
                  value

## Analysis of Songs By Taylor Swift + Her features

df %>%

Select(song, artist, rank, weekly\_rank = "weeks.on.board")  
filter(grep("featuring Taylor Swift", artist) | grep("Taylor Swift", artist))

group\_by(song) %>%

summarise(highest\_week\_rank = max(weekly\_rank))

arrange(desc(highest\_week\_rank))

X — X — X —

## # More Advanced

date and Time

\* Data cleaning → Lubridate and stringr

\* Reshaping Data → TIDYR [Data Tidying]  
↳ Cleaning/Null values

\* Joining/Merging Datasets →

## Lubridate

Parse Date-Time (Convert strings, char, no, to Date-Times)

ymd(col) | Date | Time | y-m-d | ymd-hms()  
↓ y-m-d ↓ h-m-s ↓ y-m-d | ymd-hms()

To extract | → quarter Day  
date(n), year(\*), day(x), wday(1), qday(\*)

hour(x), minute(n), day of week

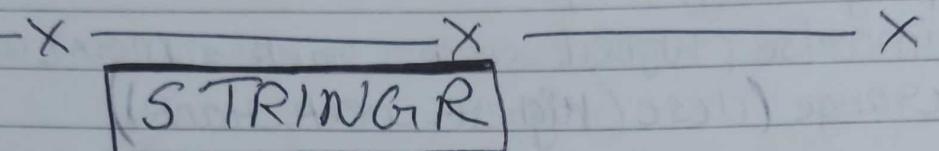
seconds(x)

week(x), quarter(n)

## Round Date-time

floor-date(x, 'month'/'second'/'year')

└ round-date | ceiling-date



Detect Matches → Str-detect(string, pattern)

→ Str-which(string, pattern) → find indexes of strings that contain a patternmatch.

→ Str-Count(string, pattern) → Count the no. of matches in a string.

## Subset Strings

→ Str-Sub(string, start, end) → Extract substring from a character vector

→ Str-extract(string, pattern) → Returns the first pattern match found in each string, as a vector. → Str-Extract-all to return every pattern

→ Str-match(string, pattern) → Returns the first pattern match found in each string, as a matrix with a col for each () group in pattern.

└ Returns the Result

## • Manage Lengths

str\_length(string)

str-trunc (string, width, side = c("right", "left", "Both"),  
pad = "") → Truncate width of strings, replacing  
content with ellipsis

str-trim (string, side = c("right", "left", "Both")) →  
Trim whitespace from the start and/or end.

## • Mutate Strings

• str-to-lower → str-to-upper → str-to-title → str-replace  
(string, replacement)

Tidyr → Data Tidiness

→ mutate(PrimaryArtist = ifelse(str\_detect(artist, 'featuring'),  
str\_match(artist, "(.\*)" ))\$featuring)

← [1,2],

.\* (one or more characters) artist]

\\s (space)

[1,2] → first pt will return finding the whole match string

.\* \\s featuring

and then will return and find-return the  
required (.\*). so [1,2]

→ Tidyr

`Supply(data, function(x) is.na(x))` → `data[, 1:7] Supply(function(x) is.na(x))`

↳ To check null values in each columns.

`is.na()` → To check null values  
`sum(is.na())` → Total null values

## Tidyr

### • Handle Missing Values

- 1) `drop_na()` → Drop rows containing NA's in column
- 2) `fill(data..., direction = "l")` → fill data in NA's columns  
 (using the next or previous down value)

→ 

|   |    |
|---|----|
| 1 | NA |
| 2 | NA |
| 3 | NA |

 → 

|   |   |
|---|---|
| 1 | 1 |
| 2 | 1 |
| 3 | 3 |

 | Fill (down)

- 3) `replace_na(data, replace)` → replace with max, mean, median or mode.

### • Functions on cells (rows).

- 1) `unite` → `unite(data, col..., sep = "", remove = TRUE, na.rm = FALSE)`

↳ Collapse cells across several columns into a single column

| col1 | col2 |
|------|------|
| 19   | 99   |
| 20   | 21   |

→ 

|      |
|------|
| col1 |
| 1999 |
| 2021 |

- 2) `separate` → `separate(data, col, into, sep =, remove = TRUE)`

↳ Separate each column into several col  
 Opposite of unite

# Reshape Data

- `pivot_longer` → `pivot_longer(data, cols = a:b, names_to = , values_to = , values_drop_na = FALSE)`  
<sup>p cols new name</sup>

↳ Collapses several columns into rows, hence lengthening the Data

|      |      |
|------|------|
| 1999 | 2000 |
| Date | Data |

`pivot_longer(data, cols = 1:2, names_to = year, values_to = cases)`

- `pivot_wider` → Inverse of `pivot_longer`

↳ widen data by expanding rows into columns

↳ `pivot_wider(data, names_from = "", value_from = "")`  
<sup>, value to impute</sup>  
↳ Col to make wider

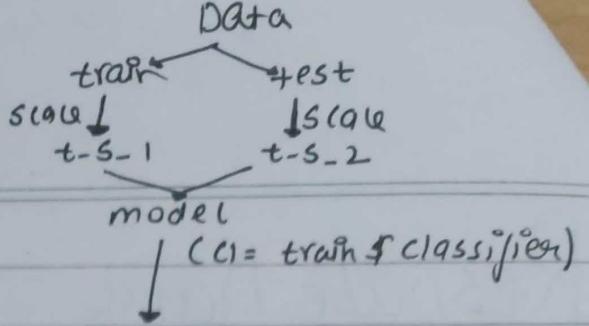
↳ `pivot_wider(data, names_from = "Type", values_from = "count")`.

Merging data frames  
X                    X                    X  
Inner Join      Left Join      Right Join      Full Outer Join



Inner\_join(DataSet2, by = c("colnamefrom01") = "colfrom02")

## KNN



### library(class)

library('caTools') → for splitting data

evaluate using [test & classifier]

data <- iris

# Data split (sample.split → Obj → use in subset())  
set.seed(123)

Split <- sample.split(data, splitRatio = 0.8) // Obj  
train-data <- subset(data, split == 'TRUE')  
test-data <- subset(data, split == 'FALSE')

### # Feature scaling (scale)

train-data-scale  
train-data <- scale(train-data[, 1:4])  
test-data <- scale(test-data[, 1:4])  
test-data-scale

### # Model

Model ⇒ Knn (train = train-data-scale,  
test = test-data-scale,  
cl = train-data\$species,  
K = 5)

conf-matrix = table(test-data\$species, model)  
conf-matrix

Acc = 1 - mean(model != test-data\$species)  
Acc

hQ.omit()

removes rows with null values

## K-means

library (factoentra)

## Library Cluster

```
df <- USArrests
```

af

```
df <- na.omit(df) # Removes Rows with null values
```

`df <- scale(df) # scaling`

df

# Elbow curve to find optimal value of K (viz - nbclust)

```
fviz_nbclust(dt, kmeans, method = "gap-start") # gap-statistic
```

`fviz_nbclust(df, kmeans, method = "wss")` // sum of square

# Model with K=4

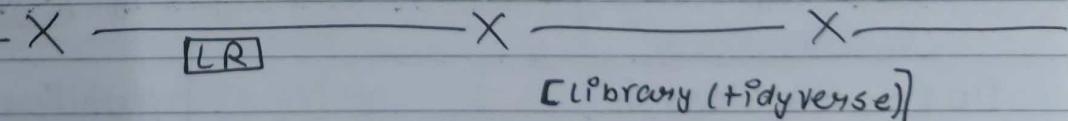
Set.seed(123)

```
model <- Kmeans(dt, centers=4, nstart=25)
```

model

## # Graph

`fviz_cluster(model, data = dt)`



Model  $\leftarrow \text{Im}(\text{y} \sim \text{x}, \text{data} = \text{df})$

`summary(model)` — R-score and all other summary

`coef(model)` → Intercept value and n ( $y = mn + b$ )  
(1) (2)

```
ggplot(data=dt, aes(x=x, y=y)) +
```

geom-point () +

gewöhnliche (intercept = (coef(modell))[1], slope = (coef(modell)[2])).

everything is  
an object

**OOPS IN RT**

- Modular pieces of code → build blocks for large systems.
- Classes - Objects - Abstraction - Encapsulation - Inheritance - Polymorphism.

**53** → overload any func<sup>n</sup> [No Inherit or encapsulation  
No safety]

**54** → Important characteristic of OOPS. [Build after 53  
Add safety]

Nowadays we use **R6**

Classes: user-defined data types that serve Blueprint for creating  
Objects - attributes - methods

Object → Instances of individual classes

- Constructor  $\Rightarrow$  Initialize = function ()
- public = list()
- private = list()
- Inherit = Parent-class name

↓  
use Blueprint to create  
separate instance with  
different values.

Inside the constructor → super \$ Initialize (s1, s2, s3) // calling Base  
class constructor

- Personal Info = Private → access them by using private \$ to access them
- Personal Info = Public → access them by using ~~public~~ **self \$**

Dog ← R6 Class C

classname = " ",

private = list()  
name = NULL,  
age = NULL,  
hair-color = NULL), →

public = list()  
Initialize = function (name = NA,  
age = NA, hair-color = NA) {

private \$ name = name  
private \$ age = age  
private \$ hair-color = hair-  
color,

bark = function()

d

y cat(private\$name, ' says fuck off')

)

)

## || Obj creation

Dog1 <- Dog\$new('Milo', 4, 'white').

Dog1\$bark()

Dog1\$name → NULL as it is Private.



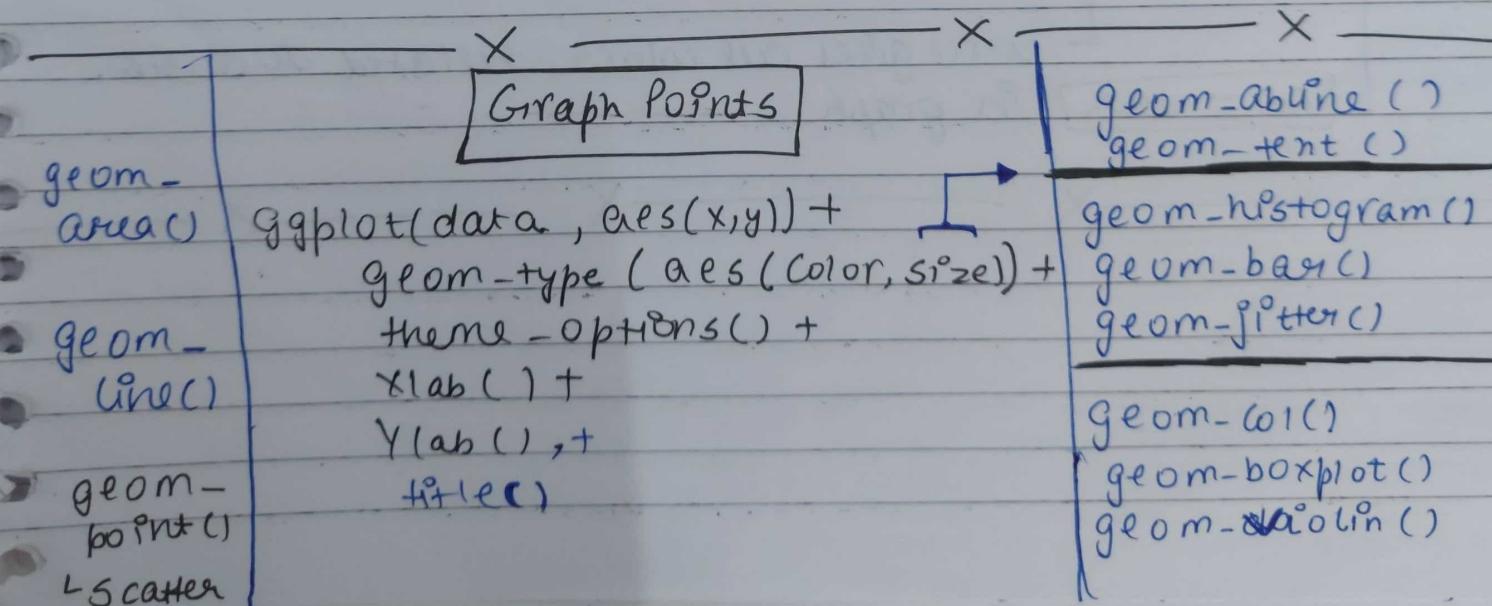
In Inheritance create a Parent class → A

In its child B Add → Inherit = A.

In B's public constructor Call parent constructor  
(super \$initializer)

will initialize the properties ←

Inherited from Parent.



## Normal

- Bar chart

barplot(data, xlab, ylab, main, col = , border = )  
 ↗ Bar color  
 ↗ title

- Histogram

hist(data, xlab, ylab, main, col = , border = ,  
 breaks = )  
 ↗ Data distribution over graph

- Box plot

boxplot(y ~ x, data, xlab, ylab)

- Scatter plot()

plot(Speed ~ dist, data, xlab, ylab, main, pch = ,  
 cex = , col = )  
 ↗ Point size  
 ↗ design  
 ↗ of points  
 ↗ (., o, △, □)

[palette()] → use for colors [or] [colors()] [has more colors]

var <- grep('dark', colors(), value = TRUE)

↳ will give all colors related to dark.  
 ↗ col = var in graph

X — X — X —