 	[<	[	٤	E	j. '	h
Ho		-	٠	٠		•

- How to work with data in R?

- Stored in R [data frame]

- Read in data into R

- Simulate samples from a given distribution.

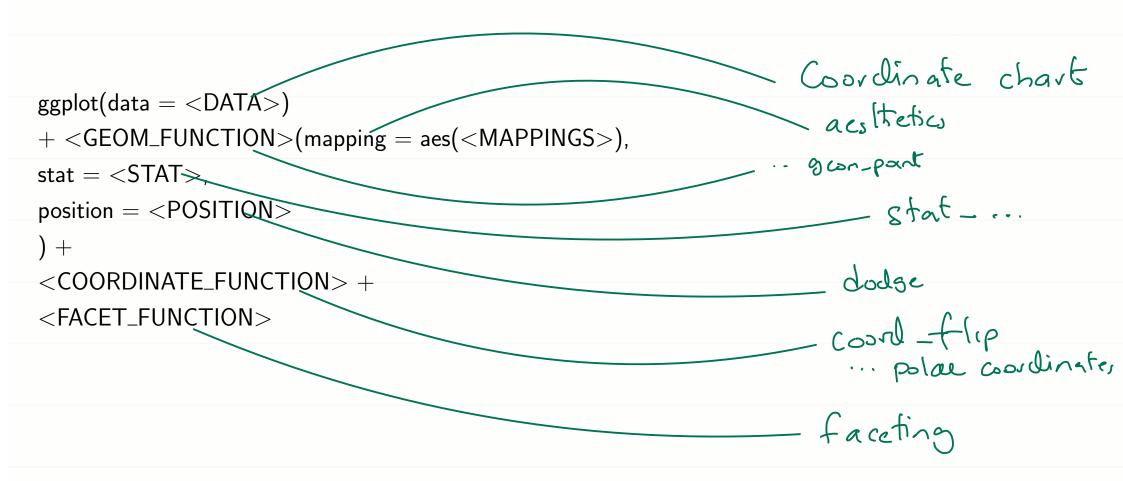
Recall - Weck 1 and 2

- How to work with R and R studio?

- Data Visualisation:
- goplet

plat

# Each template takes 7 statements/ Parameters



### Data Types in R

```
> Course = "B.Sc."
> Number = 40
> Smart = TRUE
> mode(Course)
[1] "character"
> mode(Number)
[1] "numeric"
> mode(Smart)
[1] "logical"
```

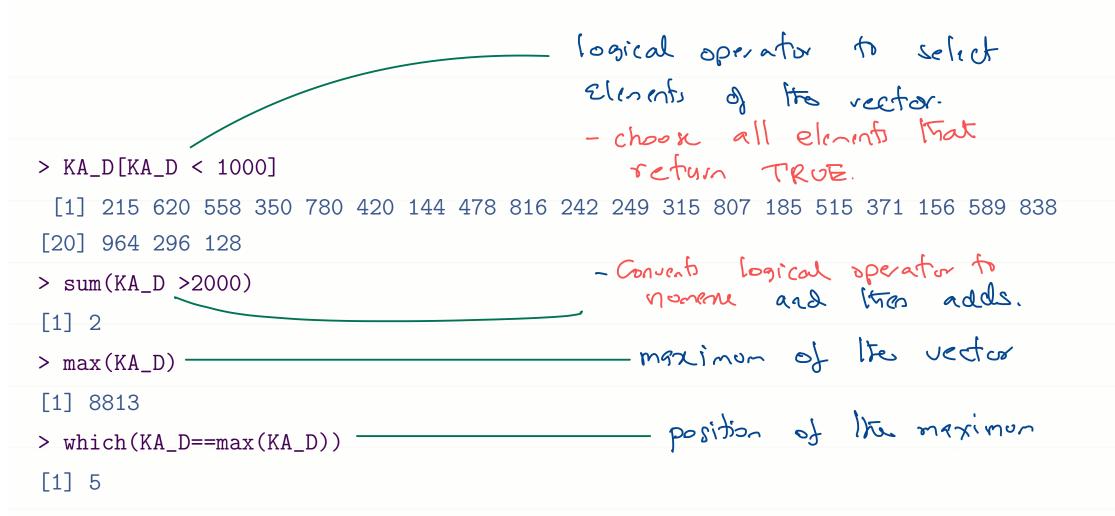
```
R has many data types - Focus on three
```

- · Character data = Surrounded by double quotes
- · logical data = TRUE or FALSE
- · Another important characteristic is class. For this we can use the class function.

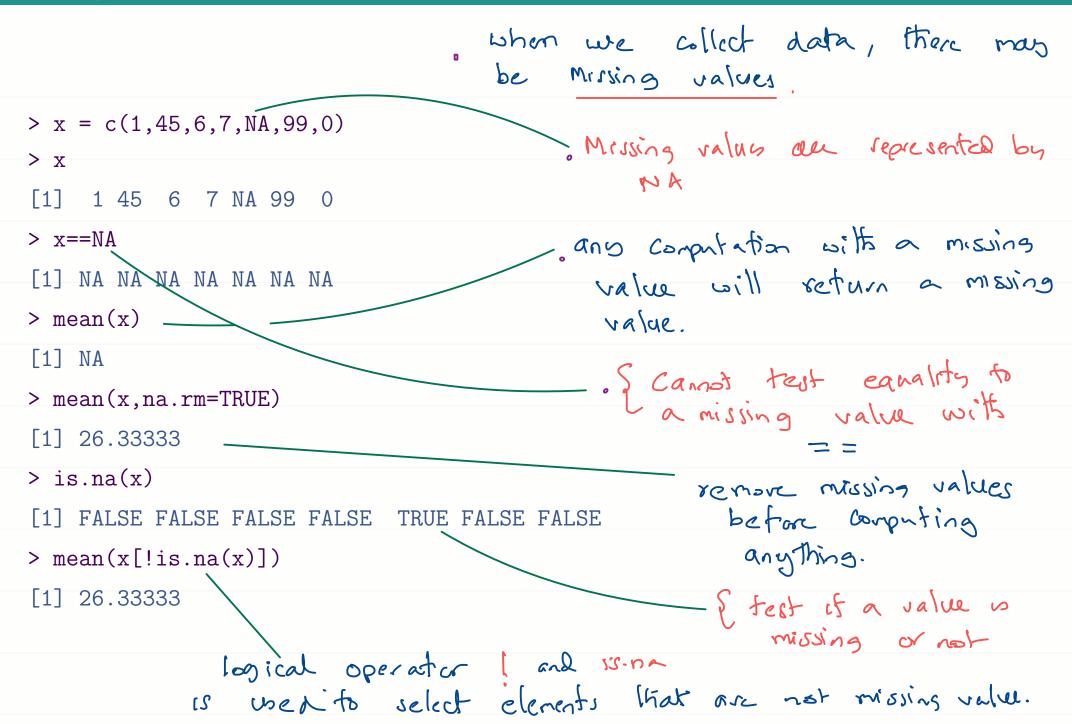
- Useful to have variables in R > x = 3:7to recall the values is vouiable > x check help [1] 3 4 5 6 7 - create a rector using sear(...) options > s = seq(1,10, by=1)> s [1] 1 2 3 4 5 6 7 8 9 10 - Position of 8.5 is 16 in > s10 = seq(1,10,by=0.5)> s10the vector [1] 10 1.5 2.0 2.5 3.0 3.5 4.0 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0 [16] <u>8.5</u> 9.0 9.5 10.0 repeats the specified element > rep(6,7)in the first argument 6,00 [1] 6 6 6 6 6 6 6 > rep(x,3)many times as the second [1] 3 4 5 6 7 3 4 5 6 7 3 4 5 6 7 argument 7. first argument can be a vector.

```
vector.
                                                            Districtuse
                                                            Discharge data of
> KA_D = c(215,620,558, 1109,8813,350, 780, 420,
                                                            COVID-19 patients
                   144,478,816,242,1051,249,1238, 315,
+
                   807, 185, 1993, 515, 1997, 2886, 371, 156,
                                                              Icain ataka
                   589,1746,838,964,296,128)
                                                              on 81-01-202L
> KA_D[c(1,3,5)]
                                                Select elements 1,3,5
        558 8813
     215
                                                 Renaval elentro 1:20
> KA_D[-c(1:20)]
 [1] 1997 2886 371
                     156
                         589 1746
                                   838
                                         964
                                              296
                                                  128
> KA_Dp = (KA_D)/sum(KA_D)*100
                                                   - vector of percentage
               rectorigation computation.
                                                    et discharge acrow distracts.
> KA_Dp
 [1]
                            1.8076387
      0.6964916
                2.0084875
                                       3.5926010 28.5496777
                                                             1.1338236
 [7]
      2.5268068
                1.3605883
                            0.4664874
                                       1.5484791
                                                  2.6434287
                                                             0.7839580
[13]
      3.4047102
                0.8066345
                            4.0104960
                                       1.0204412
                                                  2.6142732
                                                             0.5993067
[19]
     6.4563154
                 1.6683404
                            6.4692734
                                       9.3491853
                                                  1.2018530
                                                             0.5053614
[25]
      1.9080631
                 5.6561599
                            2.7146976
                                       3.1228741
                                                  0.9588908
                                                             0.4146555
```

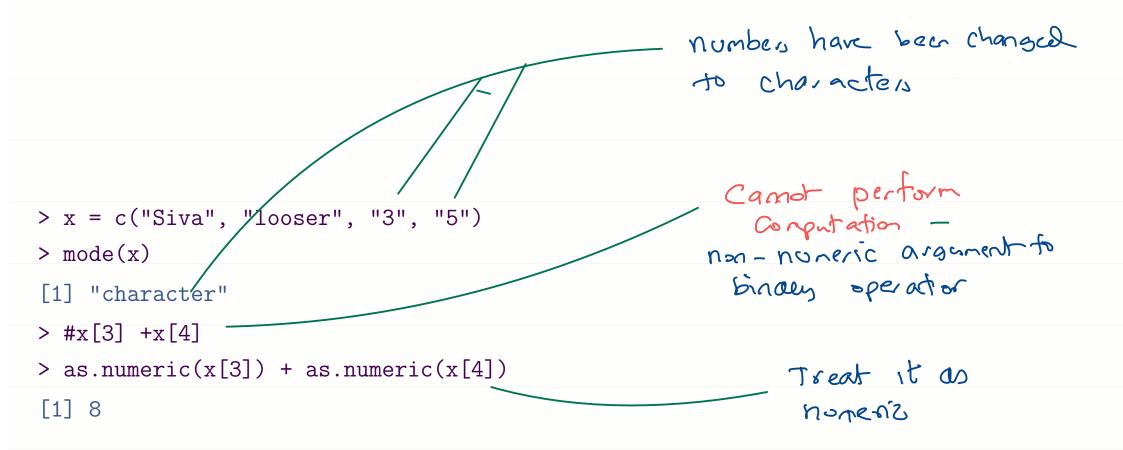
#### Vectors in R



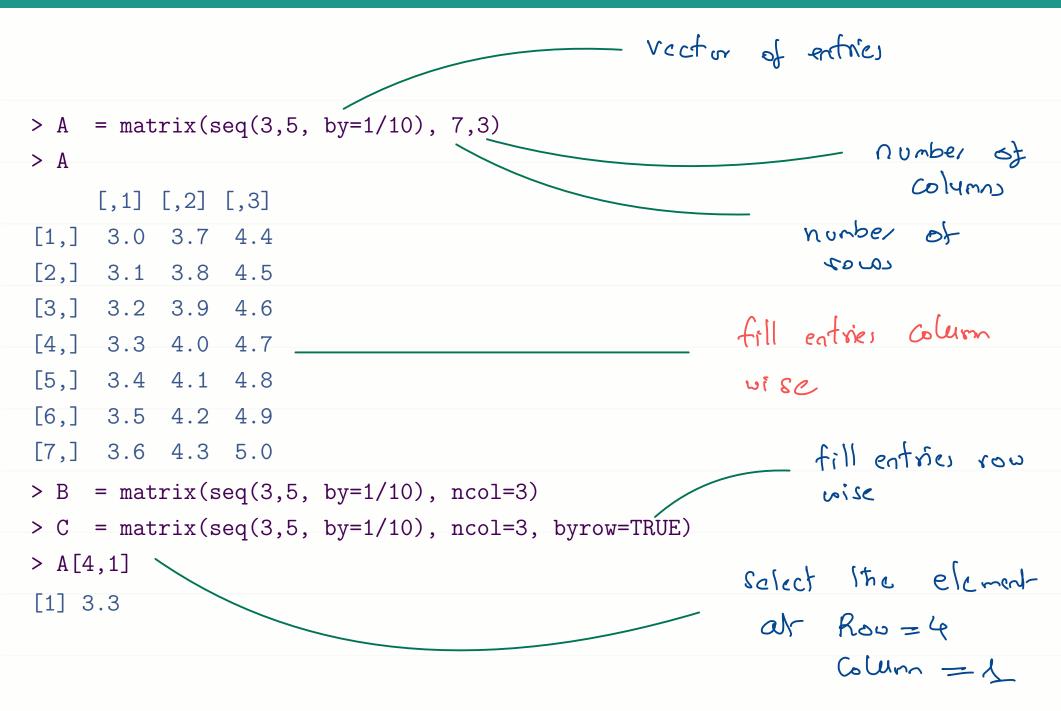
## Missing Values in R



#### Modes in R



### Matrices in R = all entries have to be in the same mall



#### Data Frames in R

- Data frame is like a matrix I rectangulær - Each column in a data frame Asrany can be in a different mode

```
> xd = c("Siva", "looser", 3, 5)
> xd

[1] "Siva" "looser" "3" "5"

be in the same made
```

Datafrance - best was to store data in R

# Data Frames in R - our first creation

			١					
	lated a			(coore sponding)				
Mac.	mes of	distric	£1.					
> KA_District=c("Bagalakote", "Ba	llari","Bela	agavi","E	Bengalur	ru Rural","Bengaluru Urban"				
+ "Bidar", "Chamarajanagara", "Chikkaballapura", "Chikkamagaluru", "Chitradurga",								
+ "Dakshina Kannada", "Davanagere", "Dharwada", "Gadag", "Hassana", "Haveri",								
+ "Kalaburagi", "Kodagu", "Kolara", "Koppala", "Mandya", "Mysuru", "Raichuru",								
+ "Ramanagara", "Shivamogga", "Tumakuru", "Udupi", "Uttara Kannada", "Vijayapura"								
+ ,"Yadagiri"			~ L					
+ )			. Creat	es a data frame				
> KA_Discharge = data.frame(KA_District, KA_D) Discharge data.								
<pre>&gt; class(KA_Discharge)</pre>			- Didx	nd names				
[1] "data.frame"								
> mode (KA_Discharge)  Secution on data frace								
[1] "list"	C> 10 1c			applies like				
> sapply(KA_Discharge,mode) <	Simpli fie	tractico	6 g d	male function				
KA_District KA_D	120p	7 321 47						
"character" "numeric"				to each				
vousable (collisson 19								
			dato	- Kiden (-				

## Data Frames as Matrix in R - working with them.

```
the names
                                                      the voerables in
                                                           data tranc
> names(KA_Discharge)=c("District", "Recovered")
> KA_Discharge$Recovered
 [1]
          620
               558 1109 8813
                             350
                                  780
                                       420
                                             144
                                                 478
                                                      816
                                                           242 1051
                                                                     249 1238
[16]
     315
          807
               185 1993
                        515 1997 2886 371
                                             156
                                                 589 1746
                                                           838
                                                                964
                                                                     296
> KA_Discharge[3,2]
                                       Select :- objects, Rows and
[1] 558
                                             Columns from Data France
> KA_Discharge[3,]
 District Recovered
                558
3 Belagavi
> KA_Discharge[,"Recovered"]
          620
               558 1109 8813
                              350
                                   780
                                        420
                                             144
                                                 478
                                                      816
                                                           242 1051
[16]
               185 1993
                         515 1997 2886
                                       371
                                                  589 1746
     315
          807
                                             156
                                                           838
                                                                964
                                                                     296
                                                                          128
```

## Data Frames as Matrix in R - add another voerable to data frame

```
> Deaths= c(346, 1712, 975, 903, 16593, 407,515, 446, 400,221, 1750, 611,
   1333,328,1291,652,856, 343, 647, 530, 673, 2494, 346,
   338, 1105, 1172, 509, 793, 500, 206
+ )
> KA_Discharge$Deaths = Deaths <
> head(KA_Discharge)
        District Recovered Deaths
     Bagalakote 215
                            346
1
        Ballari 620
                           1712
     Belagavi 558 975
4 Bengaluru Rural
                    1109 903
5 Bengaluru Urban
                    8813
                          16593
6
          Bidar
                     350
                            407
```

· Creates a rociable "Dealths" in Theo

data frame KA-Discharge fill the entre Doalbs.