

## UNIT-III

## DATA FRAMES

Data frames :-

A Data frame is like a matrix with a two-dimensional rows & columns structure. It differs from a matrix in that each column may have a different mode.

A data frame is a list with component of that list being equal length vectors. R does allows the component to be other type of objects including other data frames. In this all components of data frame are vectors.

Creating data frames :-

```
eg:- > kids <- c("x", "y")
      > Ages <- c(11, 12)
      > d <- data.frame(kids, ages, stringsAsFactors = FALSE)
      > d
      0.0 0.0
      0.1 0.1
      1.0 1.0
      o/p:   kids  age
              x    11
              y    12
```

The first two arguments in the call to `data.frame()` is clear. The third argument `stringsAsFactors = FALSE` requires, because if the named argument `stringsAsFactors` is not specified then by default `stringsAsFactors` will be true. This means that if we create a data frame from a character vector `kids` R will convert that vector to a factor. That's why here we will said `stringsAsFactors = FALSE`.

Accessing data frames:

We can access it such as via component index values or component names

eg:- > d[[]]

o/p: "x" "y" [1] x y

> d\$kids

o/p: "x" "y" [1] x y

> d[, 1]

o/p: "x" "y" [1] x y

> str(d)

data.frame 2 absolute of 2 variables

  kids: string "x" "y" 0.0

  age: number 11 12

## read.table

```
"Enam1" "Enam2" "Qwe"
```

2.0	3.3	4.0
3.3	2.0	3.7
4.0	4.0	4.0
2.3	0.0	3.3
2.3	1.0	3.3
3.3	3.7	4.0

> examplequiz <- read.table ("Enam4", header = TRUE)

> head (examplequiz)

O/p: "Enam1" "Enam2"

2.0	3.3	(21,1) → 200b <
3.3	2.0	2.0
4.0	4.0	4.0
2.3	0.0	b <
2.3	1.0	1.0
3.3	3.7	3.7

## Other matrix like operations:

Various matrix operations also apply to data frames. We can do filtering to extract various sub data frames.

Extracting sub data frames: A data frame can be viewed in rows and columns terms we

can extract sub data frame by rows, cols.

> examsquiz [2:5,]

O/p: 3.3 2.0 3.7

4.0 4.0 4.0

2.3 0.0 3.3

2.3 1.0 3.3

> Enamsquiz [2:5,2]

O/p: 3.3 2.0

4.0 4.0

2.3 0.0

2.3 1.0

> class (examsquiz [2:5,2]) O/p: numeric

> Enamsquiz [2:5, 2, drop = FALSE]

1 2.0

2 4.0 To convert to smart.Bstab

3 0.0 "Y" "X" smart: 2014

4 1.0 "Y" "X" smart: 2014

> class (Enamsquiz [2:5, 2, drop = FALSE]) O/p: data.frame

(b) if not Y takes vector  
H2 < formate.

## More on Treatment of NA values: 26/04/2022 (54)

Suppose the 2nd exam score for the first student had been missing 2.0 NA 4.0  
in any subsequent statistical analysis R would do its best with the missing data `na.rm=TRUE` explicitly telling R to ignore NA values.

Eg: 2.0 NA 4.0

> n <- c(2, NA, 4)

> mean(n)

[1] NA

> mean(n, na.rm=TRUE)

[1] 3

> examsquiz[examsquiz\$Exam.1 >= 3.8,]

o/p: 4+0

> subset(examsquiz, Exam.1 >= 3.8)

o/p: 4+0

- searched in matrix form

- using rows.cols

- function is different

The columns names are taken in the context of the given data frame the data frame of any observation that has atleast one NA value a function for this `complete.cases()`

> d+

kids states

Jack T

NA AP

Jhon MP

Suman <NA>

If we take empty row then  
it will point NULL

> complete.cases(d+)

[1] TRUE FALSE TRUE FALSE

Using the `rbind()` and `cbind()`

In this using `rbind()` to add a row

(d+)

kids ages

Jack 10

Jhon 12

NA

> rbind(d+, list("Lavanya", 19))

> Eq <- cbind(examsquiz, examsquiz\$Exam.2)

- examsquiz\$Exam.1

-1.3

0

-2.3

-1.3

0.4

> Eq

> examsquiz & Examsdiff ← (examsquiz & examsqa - examsquiz  
\$exam.1) / 2

> head(examsquiz)

o/p:	Examsquiz	Examsqa	Examsqa - examsquiz	Exam.1
1.3	2.0	3.3	4.0	3.7
-1.3	3.3	4.0	4.0	3.3
0	4.0	0.0	3.3	3.3
-2.3	2.3	4.0	3.3	3.3
-1.3	2.3	3.7	4.0	3.3
0.4	2.3			

24/04/2022

>d  
kids ages (first adding column without using bind)  
Jhon 12  
Jack 10  
>d\$one ← 1

o/p: kids ages one  
Jhon 12 1  
Jack 10 1

Applying the apply function:

> apply(examsquiz, 1, mean)

o/p: 4.0 4.0 3.7 4.0 3.3 3.3 4.0 more than 6  
3.4  
4.0  
3.3  
3.3  
4.0 gives matrix format  
T gives vector  
TA gives transpose  
TM gives matrix format  
TA > TM

Merging data frames:

In R two data frames can be similar. Combine using the merge()

merge(x,y) this merge data frames. x & y it assumes that two data frames have one or more columns with names in common.

> d1  
kids states

Jack CA  
John MA

Suman MA  
Lawanya HI

> d2  
ages kids

10 John  
12 Jack

10 Jhon  
12 Jack

> d ← merge(d1, d2)

> d

o/p: kids states ages

John MA 10  
Jack CA 12

It prints only similar

Output  
# Checks variables

p34

The two data frames have the variables kids in common & from founds the rows in which this variable had the same value of kids in both data frames.

(61)

The merge() has named arguments by.x and by.y which handles cases in which variables have similar information but different names in two data frames.

yds

ages pals

12	Jack
10	John
7	Sofia

2016-2017

11 2017

01 2017

(+0.026) (play + 3, 10)

> merge(d1, d3, by.x = "kids", by.y = "pals")

o/p: kids states ages

John	MA	10
Jack	CA	12

Even though our variables was called kids in one data frame and pals in another data frame it was mean to store the same information and does the merge made sense

Duplicate matches will appear in the full. In the result, possibly in undesirable ways.

> d2a <- bind(d2, first('15', "John"))

> d2a

o/p: kids states ages

John	10
------	----

Jack	12
------	----

Sofia	7
-------	---

John	15
------	----

> merge(d1, d2a)

o/p: kids state ages

Jack	CA	12
------	----	----

John	MA	10
------	----	----

John	MA	15
------	----	----

There are two Johns in d2a there is a John in d1 who resides in MA and another John with unkown residency in the previous merge (d1, d2a) there was only one John in both data frame. but here in the call merge (d1, d2a) it may have been the case that only one of the John was a MA. It is clear from these eg that we must use matching variable with great care

## Applying functions to data frames

28/04/2022

The `lapply` and `sapply` functions work with data frames using `lapply` and `sapply` on data frames with a specified function `f(x)` the `f(x)` will be called on each of the frame column with the return value placed in a list.

as.data.frame(d)

```
>d  
kids ages  
Jack 12  
John 10
```

John	10
Jack	12

```
>d1<-lapply(d,sort)  
>d1
```

Prints in data frame  
Postmate(madam)

kids	ages	\$kids	o/p: Jack John
John	10	\$ages	o/p: 10 12 10 12
Jack	12		10 12 10 12

Characteristics of data frames:

- The column names should be non-empty.
- The row names should be unique.
- The data frame can hold the data, which can be numeric, character or of factor type.
- Each column should contain the number of positive data items.

Factors and tables:

- Factors forms the bases for R's powerful operations including many of those performed on tabular data.
- Factors comes from notation of nominal, categorical, variables in statistics.

These values are non-numerical in nature corresponding to categories such as republican, democrat, unaffiliated etc thus they may be coded using numbers.

factors and levels

29/04/2022

An R factor might be viewed simply as a vector with a bit more information added. the extra information consist of a record of the distinct values in that vector called levels

eg: `x1 <- c(5,12,13,12)`

`xtf <- factor(x1)`

at now `print(xtf)` o/p: [5, 12, 13, 12]   
from R factors [1]

Levels 5 12 13 but 5 is first 12 is second 13 is third

It isn't 6 but 10 is first 12 is second 13 is third

of the statistics problem we have our first p3 part mark is 20

```

> str(xf) → o/p: "5" "12" "13" 1 2 3 2 (63)
(→ undass(xf) → 1 2 3 2
> attr("levels"). o/p: "5" "12" "13"
> length(x) o/p: 4
> length(xf) o/p: 3
> level(xf) o/p: 3 "5" "12" "13"

```

The base of xf here is not (5, 12, 13, 12) but rather 1 2 3 2. The letter means that our dat consist first of a level 1 value and level 2 and level 3 values and finally another level 2 value. So, the data has been recorded by level. The levels them selves are recorded to the length of the factor is still defined in terms of the length of the data rather than being a count of the no. of the levels.

```

> x <- c(5, 12, 13, 12)
> xf <- factor(x, levels = c(5, 12, 13, 12))
> nff
(→ print(o/p: 5 12 13 2)
* (original) print(5 12 13 12)
* (actual) print(5 12 13 12)
> nff[2] <- 88
> nff
o/p: 5 88 13 12

```

⇒  $x \leftarrow c(1, 2, 3, 4, 1)$  xf <- factor(x)  
> length(x) o/p: 5  
> factor(length(x)) o/p: 5  
> level(x) o/p: 4

factor(x) o/p: {1} 1 2 3 4 5

levels : 1 2 3 4

print(levels(factor(x))) = o/p: {4}

> array(x) o/p: 1 2 3 4

> length(array(x)) o/p: 5

> xf <- factor(x, levels = c(1, 2, 3, 4, 1))

> xf o/p: {1} 1 2 3 4

levels : 1 2 3 4

(64)

```

> m <- c(5, 12, 13, 12)
> mff <- factor(m, levels = c(5, 12, 13, 88))
> mff[4] 5 12 13 88
Levels 5 12 13 88
> mff[2] <- 88
> mff[4] 88 12 13
levels 5 12 13 88

```

mff did not contain the value 88 but in defining it we allowed for that future possibility. we did add the value by the same token. If we can add in an illegal level it displays like

```

> mff[2] <- 88
Warning message:
In [<-, factor, 1:step <= 12, value = 88] :
  invalid factor level, NA generated

```

### Common functions used with factors.

Here we have know about the member of the family of apply() i.e tapply as well as two other functions commonly used with the factors split(), by()

Tapply():

```

> ages <- c(25, 26, 55, 37, 21, 42)
> affils <- c("R", "D", "D", "R", "V", "D")
> tapply(ages, affils, mean)

```

We have a vector  $x_6$  of ages of voters and a factor  $f$  showing some non-numeric values such as party affiliations. We might wish to find the mean ages in  $x_6$  within each of the party group.

To print a table:

```

> d <- data.frame(gender = c("M", "M", "F", "M",
                             "F", "F"),
                     ages = c(47, 59, 21, 32, 33, 24),
                     income = c(55K, 88K, 32K,
                               52K, 34K, 24K))

```

(65)

&gt;d

	gender	ages	Income\$
1	M	47	55K
2	M	59	88K
3	F	21	32K
4	M	32	52K
5	F	33	34K
6	F	24	24K

> d \$ over 25 ← if else (d\$age > 25, 1, 0)

&gt;d

	gender	ages	income	Over 25
1	M	47	55K	1
2	M	59	88K	0
3	F	21	32K	0
4	M	32	52K	1
5	M	33	34K	1
6	F	24	24K	0

> tapply (d\$income, lapply (d\$gender, \$over25), mean)

F 2800.00 34000.00

M NA 65000.00

Split() 30/04/2022

Split is a Vector which splits the vector into group and then applies a specified function on each group.

Split function stops at that first stage just forming the groups

S(a,f) here a and f playing roles similar to those in the call tapply(m, f) i.e. m being a vector or data frame and f being a factor or a list of factors. The action is to split m into groups which are returned in a list.

eg:- >d gender ages Income Over 25

M	47	55K	1
M	59	88K	0
F	21	32K	0
M	32	52K	1
F	33	34K	1
F	24	24K	0

> split(d\$income, list(c\$gender, d\$c over 25))

Op: F		Op: M		Op: A		Op: I		Op: P		Op: C	
34K	24K	100	100	100	100	100	100	100	100	100	100
M	O	100	100	100	100	100	100	100	100	100	100
A/A	*	100	100	100	100	100	100	100	100	100	100
F	I	100	100	100	100	100	100	100	100	100	100
34K	24K	100	100	100	100	100	100	100	100	100	100
M	I	100	100	100	100	100	100	100	100	100	100
55K	88K	52K	100	100	100	100	100	100	100	100	100
SPLIT &	Yg ← c("M", "M", "F", "F") SPLIT	100	100	100	100	100	100	100	100	100	100
	> SPLIT(1:4, 9)	100	100	100	100	100	100	100	100	100	100
\$F	100	100	100	100	100	100	100	100	100	100	100
3	100	100	100	100	100	100	100	100	100	100	100
\$M	100	100	100	100	100	100	100	100	100	100	100
1 2 4	100	100	100	100	100	100	100	100	100	100	100
(NB: m, l, r, n, s, t, b, v, w, x, y, z)	100	100	100	100	100	100	100	100	100	100	100

The results shows the female cases are in records 3 the male cases are in records 1, 2, 4

Here we

text() txt()

Here we create the txt file the file contains the word=word sim times and climate was their 10<sup>3</sup> times  
Then the word & Climate will have the 2 levels of txt.  
It has 2 levels of the text  
here we call the SPLIT() then determine where these  
and the other word appear in text.

> findwords ← function(tf)  
+ bns p. not (tf) &  
+ this &  
+ txt ← scan(tf, "w")  
words ← split(1:length(tf), txt)  
return(words)

Op:

world

w 6

d 6

r 6

i 16

d 6

e

:

10

10

10

10

tf

c

1

1

1

1

pe

1

1

1

1

1

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

100

by()

02/05/2022

by() function works like tapply() function but it is applied to objects rather than vectors.

> abc <- read.table("example.csv", header = TRUE)

> by(abc, abc\$gender, function(m) .Im(m[,2] ~ m[,3]))

abc\$gender : f

(B<sub>0</sub> + B<sub>1</sub>) Exams

on

coefficient

$$[26 + 30000]_6 = 180156.11^*$$

	F	M
ages	26	46
Incomes	30000	65000

$$[46 + 65000]_6 = 390276.1^*$$

By function very similar to tapply() with the first argument specifying our data, the 2nd argument is grouping factors and the 3rd the function to be applied to each group working with tables

Exploring R tables

Example 3: f1 <- table(c(5, 12, 13, 12, 13, 5, 13))

c("a", "bc", "a", "a", "bc", "a", "a"))

> tapply(f1, length)

a	bc	"0"
5	2	NA
12	1	"2"
13	2	"1"

Here tapply() function breaks into sub vectors and then applies the length() function to each sub vector. this sub vector lengths are the count of the occurrences of each of the  $3 \times 2 = 6$  combinations of the two factors. for instance 5 occurs twice with "a" and not at all with "bc" hence the entry is 2 and NA in the first row of the output. In statistics this is called as contingency table?

- In the above example there is one problem. that is NA value. it really should be zero. meaning that in no cases with the first factor have level 5 and the second have level bc did the factor have level 5 and the second have level bc
- But the table() function creates contingency tables correctly.

> table(f1) at top also output BIB soft of

((a/b/c/freq, etc)))

5 2 bc

8 12 11 12

13 2 1

ct.csv

"Vote for x" "Voted for x last time"

"Yes"	"Yes"
"Yes"	"No"
"No"	"No"
"Not sure"	"Yes"
"No"	"No"

In the usual statistical analysis each row in these file represents one subject under study in this case, we have asked 5 people the following 2 Questions.

1. Do you plan to vote for candidate x, 05/05/2022

2. Did you vote for x in the last election

env = <-- read.table ("ct.csv", header=T)

yt

Vote for x Voted for x last time

"Yes"	"Yes"
"Yes"	"No"
"No"	"No"
"Not sure"	"Yes"
"No"	"No"

> ettab <- table(ct)

> ettab

Vote for x Voted for x last time

No	No	Yes
Yes	2	NA
Not sure	1	1

The 2 in the upper-left corner of the table shows that we had, for example, two people who said No to the first and second question. The 1 in the middle-right indicates that 1 person not sure to first question and Yes to the and Question. We can get 10 counts, which are counts on a single factor >table (c(5,12,13,12,8,5))

5	12	13	8
2	2	1	1

An example of 3D table involving voters genders, race (white, black, Asian and others) and political views (liberal or conservative)

(69)

ex: > vt # dataframe

gender	race	pol	value
M	W	L	1
M	W	C	1
F	A	C	1
M	O	CL	1
F	B	L	1
F	B	C	1

> vt <- table(v)

> vt

pol	gender	race	value
L	F	A,B,W,O	1
L	M	O	1
C	F	O	2

dim(vt) [1] 3 2 2 dim(vt[,1]) [1] 2

R programming prints out a three dimensional table as a series of two dimensional table. In these case it generates a table of gender and race for conservatives and then a corresponding table for liberals.

For example the second 2D table says that there were two white male liberals.

Matrix or array:

Matrix or array operations can be on dataframe, they can applied on table too. for ex: we can access the table cell counts using matrix notation apply these to our voting example from the previous

Ex: > class(cttab) o/p: "Table"

> ctab[1,1] o/p: 2

> ctab[,1] o/p: 20

In the second command, even though the first command had shown that ctab had class has class (cttab) we treated it as a matrix and pointed as [1,1] element. Continuing the third command pointed the first column of these matrix.

> `ctab /s`

		Voted for X last time		
		No	Yes	
Vote for X	Yes	0.2	0.2	76
	No	0.4	0.0	16
	Not sure	0	0.2	18

we can multiple the matrix by a scalar

e.g. apply (ctab, 1, sum)

Yes	Not sure	No
2	1	2

> add margins (ctab\$)

Vote for X      Voted for X last time

No	No	Yes	Sum	
Yes	0	0	2	100
Not sure	1	1	2	100
Sum	3	2	5	100

The marginal value of a variable are those obtained when this variable is held constant while others are summed. In the voting example the marginal value of vote for X variable are:  $2+0=2$ ,  $0+1=1$  &  $1+1=2$ . Here we can obtained these via the matrix apply() function. But R's applies a function add margins () from these to find marginal totals. It shows in the above example.

> dimnames (ctab)

vote for X

"No" "Yes" "Not sure"

voted for X last time

"No" "Yes"

Extracting a sub table

containing with voting example

vote for X      Voted for X last time

No	No	Yes
Yes	2	0

Not sure	No	Yes
0	0	0

Suppose, we wish to present these data at a meeting. Concentrating on those response who know they will vote for X in the current election's i.e. we wish to eliminate

the not sure entries and present a sub-table.

> Subtable 'tbltab' Note for  $x \in c("NO", "YES", \text{Voted})$

o/p: last-time = c("NO", "YES")

o/p: Note for x Noted for x last-time

	Yes	No
Yes		
No		

The function subtable performs subtable extraction

it has two arguments.

1.  $\text{tbl} \rightarrow$  the table of interest of class "Table"

2. subnames - A list specifying the desired subtable extraction, each component of list is named after some dimensions of  $\text{tbl}$ , and value of that component is a vector of the names of the desired levels.

From the above example, the argument  $\text{tbltab}$  will be a 2D table with dimension names Note for x & Noted for x last-time within these 2D the level names are NO, Notsure, YES in the 1st dimension and NO, YES in the second dimension.

we wish to exclude the not sure cases so our actual argument corresponding to the formal argument subnames

Ex:-  $\text{subtable} \leftarrow \text{function} (\text{tbl}, \text{subnames})$

{

$\text{tblarray} \leftarrow \text{unclass} (\text{tbl})$

$\text{dargs} \leftarrow \text{list} (\text{tblarray})$

$\text{ndims} \leftarrow \text{length} (\text{subnames})$

for  $i \in 1:n\text{dims}$  {

$\text{dargs}[i+1] \leftarrow \text{subnames}[i]$

}

$\text{subarray} \leftarrow \text{do.call} ("r", \text{dargs})$

$\text{dms} \leftarrow \text{lapply} (\text{subnames}, \text{length})$

$\text{Subtbl} \leftarrow \text{array} (\text{subarray}, \text{dms}, \text{dimnames} = \text{subnames})$

$\text{class} (\text{Subtbl}) \leftarrow \text{"Table"}$

$\text{return} (\text{Subtbl})$

get

get

get

The function `table` we found that it is an array whose elements are the cell counts so the strategy is to extract the desired subarray then add names to the dimensions of the subarray and then class status to the result.

For the code here, then the first task is to form the subarray corresponding to the user's desired sub-table, and this constitutes most of the code. In line three, we first extract the full cell counts array, sorting it in `tbl` `tblarray`.

Finding the largest cells in a table:

`Tabdom()`: `Tabdom()` function reports the dominant frequencies in a table.

Syntax: `tabdom(tbl, k)`

This reports the cells in the table `tbl` that have the `k` largest frequencies

```
> d <- c(5, 12, 13, 14, 3, 28, 12, 12, 9, 5, 5, 13, 5, 4, 12)
> dtab <- table(d) → 0/p: 13 4 5 9 12 13 28
> tabdom(dtab, 3). → 1 2 4 1 4 2 1
0/p: d freq. 3 → 0/p: 5 12 13 4 3 28 9
3 5 4 → 4 4 2 2 1 1 1
5 12 (replo("9")) UED: ob → y6M6dU2
2 4 2 → (replo("2")) piggat → effit
> as.dataframe(cttab) → tabdom(cttab, 2)
```

Note for n

voted for X time

(cttab\$ctid) =

- 1 No
- 2 Not Sure
- 3 Yes
- 4 No
- 5 Not Sure
- 6 Yes

- No
- No
- No
- Yes
- Yes
- Yes

0/p: -	vote for
0	No
1	Yes
0	voted for X has
1	No
1	No
2	freq
1	2

The function tells us that the values 5 & 12 were the most frequent in d with 4 instances each, and the next most frequent value was 4 with 2 instances. Consider our table cbtab. In the example we use, tabdom (above example) (73)

so the combination "No, No" was most frequent with 2 instances which the 2nd most frequent being yes-no which is 2 instance

This is not the original data frame ct from which the table cbtab was constructed. It is simply a different presentation of the table itself. They spreads 1 row for each combination of the fact with a frequency column. Added to show the no. of instances of each combination

```
function (tbl, k)
{
  tbldf <- as.data.frame(tbl)
  freqord <- order(tbldf$freq, decreasing = TRUE)
  dom <- tail(freqord, k)[1:k]
  return(dom)
}
```

09/05/2022

- The first line finds the cells in table tbl with the k highest frequencies.
- Second line indicates the enters into the function
- 3rd line indicates the create a data frame representation of tbl, adding a frequency column
- 4th line determine the proper positions of the frequencies in a sorted order
- 5th line is & rearranging the data frame in that order and that takes the 1st k rows.
- 6th line returns the result stored in dom
- 7th line close the working functions.

Other factors and table related functions :-

R includes a no. of other functions that are handy for working with tables and factors. Here we discuss two of them - 1. is aggregate()  
2. is cut()

Aggregate()

The aggregate() function calls lapply once for each variable in a group.

- aggregate() is used to get the summary statistics of the data by group. The statistics includes sum, mean, MAX, MIN and etc.,

Syntax :- aggregate (data.frame \$ aggregate.col, list (data.frame \$ group.col), FUN)

Where, data.frame → Is the input data frame

aggregate.col → Is the column to be aggregated in the data frame

group.col → Is the column to be grouped with fun

FUN - Represents sum or mean or MAX or MIN etc.

Eg :- data ← data.frame (subjects = c("R", "DAA", "R", "SE", "DAA", "SE"), id = c(1, 2, 3, 4, 5, 6), names = c("A", "B", "C", "D", "E", "F"), marks = c(89, 89, 76, 89, 90, 67))

Point (data)

O/p :- Subject id names marks

	Subject	id	names	marks
1	R	1	A	89
2	DAA	2	B	89
3	R	3	C	76
4	SE	4	D	89
5	DAA	5	E	90
6	SE	6	F	67

> Point (aggregate (data \$ marks, list (data \$ subjects),

fun = sum))

$$- 89 + 67 = 156$$

DAA

$$- 89 + 76 = 149$$

R

$$- 79 + 76 = 155$$

Cut :-

The cut() function enables us to divide the numeric vector into a range of certain intervals in a customised fashion. With cut(), the values get divided into  $n$  intervals from the  $n$  data values depending upon the breaking interval criteria.

Syntax :- `Cut(x, b, labels = FALSE)`

Where  $x$  - Is a data i.e. the numeric vector which has to be transformed

$b$  - Indicates bins are defined to be the semi open intervals.  $\{b[1], b[2], b[3]\}$  (includes)

Labels - they are logical labels attached to the result. If `labels = FALSE`, integer code will come to use.

`x2 <- c(0.8, 0.2, 0.5, 0.4, 0.4, 0.2, 0.0, 0.8)`

`> seq(from = 0.0, to = 1.0, by = 0.1)`

O/P: 0.0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0

`> binmarks <- seq(from = 0.0, to = 1.0, by = 0.1)`

`> cut(x2, binmarks, labels = FALSE)`

O/P: [1] 9 3 6 5 5 3 1 9

The above example says that  $x2[1]$  is 0.8 fell into bin 9 which was 0, 0, 0.1.  $x2[2]$  is 0.2 fell into bin 3 and so on.

(from)  $\rightarrow$  p <  
(to)  $\rightarrow$  q <  
(by)  $\rightarrow$  r <

p <= q & r <

i  $\rightarrow$  j < r

(q  $\geq$  p) always <

p + r = q <

(q + rmq) <

if [j] &  $\rightarrow$  q <