

Assignment No. 5

- * Aim :- Perform different operations using R/python.
- * Problem statement :- Perform the following operations using R/python on the Amazon book review & facebook metrics data sets.
- * Objectives :-
 - ① To learn R/python programming
 - ② To learn different data preprocessing techniques.
- * Theory :-
 - Subsetting data :- R has powerful indexing factor feature for accessing object elements. These features can be used to select and element variables and observations. The following code snippets demonstrate ways to keep or delete variables and observations and to take random samples from a dataset.
 - selecting (keeping) variables.


```
# select variables v1, v2, v3; myvars <- {v1, 'v2', 'v3'}
newdata <- mydata [myvars]
```
 - Excluding (dropping variables)

```
# exclude variables v1, v2, v3
```

```
myvars <- names(mydata) %in% c(
  "v1", "v2", "v3")
```

```
newdata <- mydata[, myvars]
```

* selecting observations :-

first 5 observations

```
newdata <- mydata[1:5,]
```

* selection using a subset function :-

The subset() function is easiest way to select variables and observations. In the following examples we select all rows that have a value of age greater than or equal to 20 or age less than 10. We keep the ID and weight columns.

using subset function

```
newdata <- subset(mydata, age >= 20 |
  age < 10)
```

```
select = c(ID, weight)
```

- Merging data :-

To merge two data frame horizontally see the merge function. In the most cases you join two data frames by one or more common key vars.

merge two data frames by ID

total ← merge(data frame A, data frame B,
by "ID")

* sorting of data

A sorted data frame in R, use the `order()` function. By default, sorting in ASCENDING. Prepend the sorting variables by a minus sign to indicate DESCENDING order.

Ex.

`attach(mtcars)`

`newdata ← mtcars[order(mpg)]`

- Transpose

Use the `t()` function to transpose a matrix or data frame. In the latter case, rownames become variable (column) name.

Example using builtin dataset

`mtcars`

`t(mtcars)`

* The Reshape package

Basically you melt data so that each row is a unique id-variable combination. Then you cast the melted data into any shape you would like. Here is a very simple example

Example of melt function

```
library(reshape)
mdata <- melt(mydata, id=c, ("id", "time"))
```

Cast melted data

```
# cast(para, formula, function)
```

```
submeans <- cast(mdata, id-variable, mean)
```

```
time_means <- cast(mdata, time-variable, mean)
```

* Conclusion & Thus we have studied & understand different data preprocessing technique