

## ASSIGNMENT-7

AIM: Integrate R and Hadoop

PROBLEM STATEMENT: Integrate R & Hadoop & perform the following operations on forest fire data.

- Text mining in R
- Data analysis using MapReduce in R.

### THEORY:

#### 1. TEXT MINING:

Natural languages are different from programming languages. The semantic meaning of statement depends on the context, tone & a lot of other factors.

Text mining deals with helping computers understand the meaning of text. Some of the common text mining applications include sentiment analysis.

In R, packages useful in understanding text & extracting insights from text & text mining package are as follows:

- RSQLik, SQLik, for R
- tm, framework for text mining
- snowballc, text streaming library
- wordcloud for making visualisation
- lyryxhet for text sentiment analysis
- ggplot2 for data visualisation
- planteda N-grams.



## 2. TEXT PREPROCESSING :

Text data contains white spaces, punctuations, stop words, etc.

Depending on the task at hand, deal with different characters

- Convert text to lower case
- Remove numbers
- Remove English stopwords
- Remove extra whitespaces
- Eliminate punctuation marks.

## 3. CLEANING TEXT IN R :

- i. # transform & clean the text
- ii. library("tm")
- iii. doc & corpus(vector source(emailRa i))

Transformations are done with the `tm-mapf` function to all elements of the corpus.

A document term matrix is important representation for text mining in R tasks & an important concept.

## 4. WORD CLOUD :

A word cloud is a simple yet informative way to understand textual data to do analysis

```
library(wordcloud)
```

For word cloud computing terms with a frequency greater than 30, use following commands:



wordcloud (names (freq), freq, min.freq = 30,  
colours = brewer.pal (3, "dark"))

### CONCLUSION :

In this assignment, we learnt to integrate Hadoop & R & perform text mining & data analysis.