

Distributed Text Mining with tm

Stefan Theußl¹ Ingo Feinerer² Kurt Hornik¹

Department of Statistics and Mathematics,
WU Vienna University of Economics and Business¹

Institute of Information Systems, DBAI Group
Technische Universität Wien²

08.07.2009

Text Mining in R

- ▶ Highly interdisciplinary research field utilizing techniques from computer science, linguistics, and statistics
- ▶ Vast amount of textual data available in machine readable format:
 - ▶ scientific articles, abstracts, books, ...
 - ▶ memos, letters, ...
 - ▶ online forums, mailing lists, blogs, ...
- ▶ Steady increase of text mining methods (both in academia as in industry) within the last decade

Text Mining in R

- ▶ tm Package
- ▶ Tailored for
 - ▶ Plain texts, articles and papers
 - ▶ Web documents (XML, SGML, ...)
 - ▶ Surveys
- ▶ Available *transformations*: `stemDoc()`, `stripWhitespace()`, `tmTolower()`, ...
- ▶ Methods for
 - ▶ Clustering
 - ▶ Classification
 - ▶ Visualization
- ▶ Feinerer (2009) and Feinerer et al. (2008)

Motivation

- ▶ Data volumes (corpora) become bigger and bigger
- ▶ Many tasks, i.e. we produce output data via processing lots of input data
- ▶ Text mining methods are becoming more complex and hence computer intensive
- ▶ Want to make use of many CPUs
- ▶ Typically this is not easy (parallelization, synchronization, I/O, debugging, etc.)
- ▶ Need for an integrated framework
- ▶ preferably usable on large scale distributed systems

→ **Main motivation: large scale data processing**

Motivation

- ▶ Multi-processor environments and large scale compute clusters/clouds available for a reasonable price
- ▶ Integrated frameworks for parallel/distributed computing available (e.g., Hadoop)
- ▶ Thus, parallel/distributed computing is now easier than ever
- ▶ R already offers extensions to use this software (e.g. via hive, nws, Rmpi, snow, etc.)

Distributed Text Mining in R

Example: Stemming

- ▶ Erasing word suffixes to retrieve their radicals
- ▶ Reduces complexity
- ▶ Stemmers provided in packages Rstem¹ and Snowball²

Data:

- ▶ *Wizard of Oz* book series (<http://www.gutenberg.org>): 20 books, each containing 1529 – 9452 lines of text
- ▶ *Reuters-21578*: one of the most widely used test collection for text categorization research
- ▶ *New York Times Annotated Corpus*: > 1.6 million text files

¹Duncan Temple Lang (version 0.3-1 on Omegahat)

²Kurt Hornik (version 0.0-6 on CRAN)

Distributed Text Mining in R

Difficulties:

- ▶ Large data sets
- ▶ Corpus typically loaded into memory
- ▶ Operations on all elements of the corpus (so-called *transformations*)

Strategies:

- ▶ Text mining using tm and MapReduce/hive¹
- ▶ Text mining using tm and MPI/snow²

¹Stefan Theußl (version 0.1-1)

²Luke Tierney (version 0.3-3)

The MapReduce Programming Model

The MapReduce Programming Model

- ▶ Programming model inspired by functional language primitives
- ▶ Automatic parallelization and distribution
- ▶ Fault tolerance
- ▶ I/O scheduling
- ▶ Examples: document clustering, web access log analysis, search index construction, ...
- ▶ Dean and Ghemawat (2004)

Hadoop (<http://hadoop.apache.org/core/>) developed by the Apache project is an open source implementation of MapReduce.

The MapReduce Programming Model

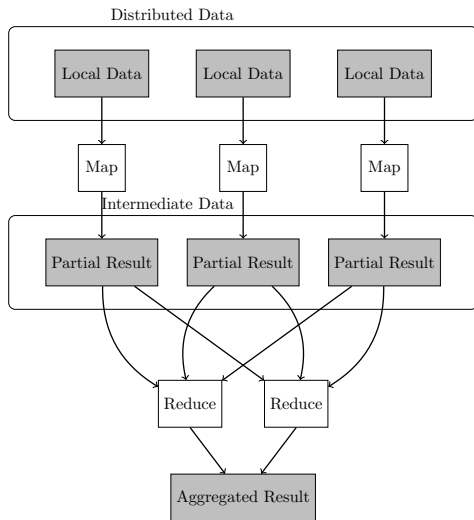


Figure: Conceptual Flow

The MapReduce Programming Model

A MapReduce implementation like Hadoop typically provides a distributed file system (DFS):

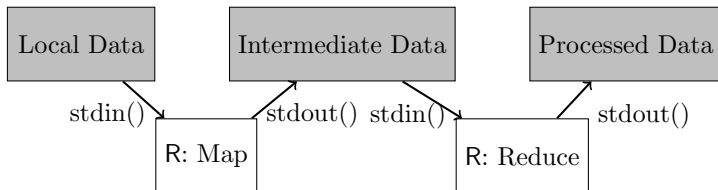
- ▶ Master/worker architecture (Namenode/Datanodes)
- ▶ Data locality
- ▶ Map tasks are applied to partitioned data
- ▶ Map tasks scheduled so that input blocks are on same machine
- ▶ Datanodes read input at local disk speed
- ▶ Data replication leads to fault tolerance
- ▶ Application does not care whether nodes are OK or not

Hadoop Streaming

- ▶ Utility allowing to create and run MapReduce jobs with any executable or script as the mapper and/or the reducer

```
$HADOOP_HOME/bin/hadoop jar $HADOOP_HOME/hadoop-streaming.jar
```

- ▶ `-input inputdir`
- ▶ `-output outputdir`
- ▶ `-mapper ./mapper`
- ▶ `-reducer ./reducer`



Hadoop InteractiVE (hive)

Hadoop InteractiVE (hive)

hive provides:

- ▶ Easy-to-use interface to Hadoop
- ▶ Currently, only Hadoop core (<http://hadoop.apache.org/core/>) supported
- ▶ High-level functions for handling Hadoop framework (`hive_start()`, `hive_create()`, `hive_is_available()`, etc.)
- ▶ DFS accessor functions in R (`DFS_put()`, `DFS_list()`, `DFS_cat()`, etc.)
- ▶ Streaming via Hadoop (`hive_stream()`)
- ▶ Available on R-Forge in project RHadoop

Example: Word Count

Data preparation:

```
1 > library("hive")
2 Loading required package: rJava
3 Loading required package: XML
4 > hive_start()
5 > hive_is_available()
6 [1] TRUE
7 > DFS_put("~/Data/Reuters/minimal", "/tmp/Reuters")
8 > DFS_list("/tmp/Reuters")
9 [1] "reut-00001.xml" "reut-00002.xml" "reut-00003.xml"
10 [4] "reut-00004.xml" "reut-00005.xml" "reut-00006.xml"
11 [7] "reut-00007.xml" "reut-00008.xml" "reut-00009.xml"
12 > head(DFS_read_lines("/tmp/Reuters/reut-00002.xml"))
13 [1] "<?xml version=\"1.0\"?>"
14 [2] "<REUTERS TOPICS=\"NO\" LEWISSPLIT=\"TRAIN\" [...]"
15 [3] " <DATE>26-FEB-1987 15:03:27.51</DATE>"
16 [4] " <TOPICS/>"
17 [5] " <PLACES>"
18 [6] " <D>usa</D>"
```

Example: Word Count

```
1 mapper <- function(){
2   mapred_write_output <- function(key, value)
3     cat(sprintf("%s\t%s\n", key, value), sep = "")

5   trim_white_space <- function(line)
6     gsub("(^ +)|( +$)", "", line)
7   split_into_words <- function(line)
8     unlist(strsplit(line, "[[:space:]]+"))

10  con <- file("stdin", open = "r")
11  while (length(line <- readLines(con, n = 1,
12    warn = FALSE)) > 0) {
13    line <- trim_white_space(line)
14    words <- split_into_words(line)
15    if(length(words))
16      mapred_write_output(words, 1)
17  }
18  close(con)
19 }
```


Example: Word Count

```
1  reducer <- function(){
2    [...]
3    env <- new.env(hash = TRUE)
4    con <- file("stdin", open = "r")
5    while (length(line <- readLines(con, n = 1,
6      warn = FALSE)) > 0) {
7      split <- split_line(line)
8      word <- split$word
9      count <- split$count
10     if(nchar(word) > 0){
11       if(exists(word, envir = env, inherits = FALSE)) {
12         oldcount <- get(word, envir = env)
13         assign(word, oldcount + count, envir = env)
14       }
15       else assign(word, count, envir = env)
16     }
17   }
18   close(con)
19   for (w in ls(env, all = TRUE))
20     cat(w, "\t", get(w, envir = env), "\n", sep = "")
21 }
```

Example: Word Count

```
1 > hive_stream(mapper = mapper,
2               reducer = reducer,
3               input = "/tmp/Reuters",
4               output = "/tmp/Reuters_out")
5 > DFS_list("/tmp/Reuters_out")
6 [1] "_logs"      "part-00000"
7 > results <- DFS_read_lines(
8       "/tmp/Reuters_out/part-00000")
9 > head(results)
10 [1] "-\t2"          "--\t7"
11 [3] ":\t1"         ".\t1"
12 [5] "0064</UNKNOWN>\t1" "0066</UNKNOWN>\t1"
13 > tmp <- strsplit(results, "\t")
14 > counts <- as.integer(unlist(lapply(tmp, function(x)
15                               x[[2]])))
16 > names(counts) <- unlist(lapply(tmp, function(x)
17                               x[[1]]))
18 > head(sort(counts, decreasing = TRUE))
19 the    to    and    of    at said
20 58    44    41    30    25    22
```

Distributed Text Mining in R

Solution (Hadoop):

- ▶ Data set copied to DFS ('DistributedCorpus')
- ▶ Only meta information about the corpus in memory
- ▶ Computational operations (*Map*) on all elements in parallel
- ▶ Work horse `tmMap()`
- ▶ Processed documents (revisions) can be retrieved on demand

Distributed Text Mining in R - Listing

```
1 > library("tm")
2 Loading required package: slam
3 > input <- "~/Data/Reuters/reuters_xml"
4 > co <- Corpus(DirSource(input), [...])
5 > co
6 A corpus with 21578 text documents
7 > print(object.size(co), units = "Mb")
8 65.5 Mb

10 > source("corpus.R")
11 > source("reader.R")
12 > dc <- DistributedCorpus(DirSource(input), [...])
13 > dc
14 A corpus with 21578 text documents
15 > dc[[1]]
16 Showers continued throughout the week in
17 [...]
18 > print(object.size(dc), units = "Mb")
19 1.9 Mb
```

Distributed Text Mining in R - Listing

Mapper, i.e. input to `hive_stream()` (called by `tmMap()`):

```
1 mapper <- function(){
2   require("tm")
3   fun <- some_tm_method
4   [...]
5   con <- file("stdin", open = "r")
6   while(length(line <- readLines(con, n = 1L,
7     warn = FALSE)) > 0) {
8     input <- split_line(line)
9     result <- fun(input$value)
10    if(length(result))
11      mapred_write_output(input$key, result)
12  }
13  close(con)
14 }
```

Distributed Text Mining in R

Infrastructure:

- ▶ Development platform: 8-core Power 6 shared memory system



IBM System p 550	
4	2-core IBM POWER6 @ 3.5 GHz
128	GB RAM

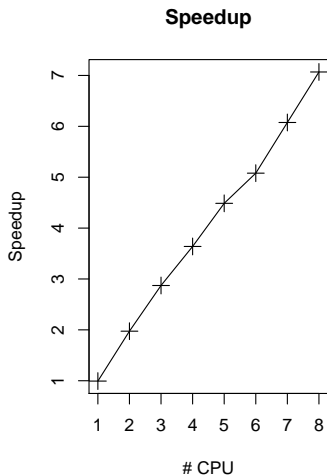
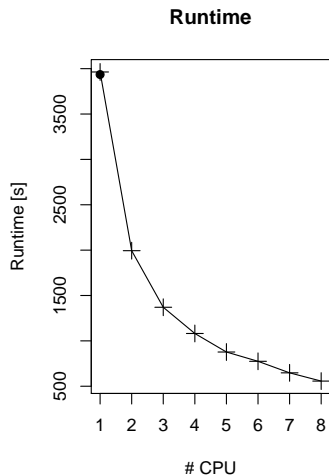
- ▶ Computers of PC Lab used as worker nodes
 - ▶ 8 PCs with an Intel Pentium 4 CPU @ 3.2 GHz and 1 GB of RAM
 - ▶ Each PC has > 20 GB reserved for DFS

MapReduce framework:

- ▶ Hadoop (implements MapReduce + DFS)
- ▶ R (2.9.0) with tm (0.4) and hive (0.1-1)
- ▶ Code implementing 'DistributedCorpus'
- ▶ Cluster installation coming soon (loose integration with SGE)

Benchmark

Wizard of Oz data set (PC Lab cluster):



Benchmark

Reuters-21578:

- ▶ Single processor runtime (`lapply()`): > 30 min.
- ▶ tm/hive on 8-core SMP (`hive_stream()`): 4 min.
- ▶ tm/snow on 8 nodes of cluster@WU (`parLapply()`): 2.13 min.

Lessons Learned

- ▶ Problem size has to be sufficiently large
- ▶ Requirement: access text documents in R via `[[`, thus location of texts in DFS important (currently: ID = file path)
- ▶ Serialization difficult: how updating text IDs? Currently via meta information to each chunk (chunk name, position in the chunk)
- ▶ Remote file operation on DFS around 2.5 sec. (significantly reduced with Java implementation)

Conclusion

- ▶ Use of Hadoop in particular the DFS enhances handling of large corpora
- ▶ Significant speedup in text mining applications
- ▶ Thus, MapReduce has proven to be a useful abstraction
- ▶ Greatly simplifies distributed computing
- ▶ Developer focus on problem
- ▶ Implementations like Hadoop deal with messy details
 - ▶ different approaches to facilitate Hadoop's infrastructure
 - ▶ language- and use case dependent

Thank You for Your Attention!

Stefan Theußl

Department of Statistics and Mathematics




email: Stefan.Theussl@wu.ac.at

URL: <http://statmath.wu.ac.at/~theussl>

WU Vienna University of Economics and Business

Augasse 2–6, A-1090 Wien

References

-  Jeffrey Dean and Sanjay Ghemawat.
MapReduce: Simplified data processing on large clusters.
In *OSDI'04, 6th Symposium on Operating Systems Design and Implementation*, pages 137–150, 2004.
-  I. Feinerer
tm: Text Mining Package, 2009
URL <http://CRAN.R-project.org/package=tm>
R package version 0.3-3
-  I. Feinerer, K. Hornik, and D. Meyer
Text mining infrastructure in R
Journal of Statistical Software, 25(5):1–54, March 2008
ISSN 1548-7660
URL <http://www.jstatsoft.org/v25/i05>