

CSCI946 Assignment

Yao Xiao
SID 2019180015

November 10, 2020

1 Task 1

```
1 library("ggplot2")
2 library("reshape2")
3 library("lda")
4
5 data(cora.documents)
6 data(cora.vocab)
7
8 theme_set(theme_bw())
9
10 K <- 10
11
12 N <- 9
13
14 result <- lda.collapsed.gibbs.sampler(cora.documents,
15                                       K,
16                                       cora.vocab,
17                                       25,
18                                       0.1,
19                                       0.1,
20                                       compute.log.likelihood=TRUE)
21
22
23 top.words <- top.topic.words(result$topics, 5, by.score=TRUE)
24
25 topic.props <- t(result$document_sums) / colSums(result$document_
26   sums)
27
28 document.samples <- sample(1:dim(topic.props)[1], N)
29
30 topic.props <- topic.props[document.samples,]
31
32 topic.props[is.na(topic.props)] <- 1 / K
33
34 colnames(topic.props) <- apply(top.words, 2, paste, collapse=" ")
```

```

35 topic.props.df <- melt(cbind(data.frame(topic.props), document=
    factor(1:N)), variable.name="topic", id.vars = "document")
36
37 qplot(topic, value*100, fill=topic, stat="identity", ylab="
    proportion_1(%)", data=topic.props.df, geom="col") +
38   theme(axis.text.x = element_text(angle=0, hjust=1, size=12)) +
    coord_flip() +
39   facet_wrap(~ document, ncol=3)

```



2 Task 2

```

1 library("dplyr")

```

```

2 library("magrittr")
3 library("data.table")
4 library("tidytext")
5 library("topicmodels")
6 library("colorspace")
7 library("purrr")
8 library("ldatuning")
9 library("gmp")
10 library("RColorBrewer")
11 library("wordcloud")
12 library("ggplot2")
13 library("lubridate")
14 library("lubridate")
15 library("reshape2")
16 library("textmineR")
17
18 data <- fread("/Users/december/Desktop/Week8/Sentiment.csv")
19
20 data <- data %>% select(text, id) %>% head(5000)
21
22 # text cleaning
23 data$text <- sub("RT.*:", "", data$text)
24 data$text <- sub("@.*_", "", data$text)
25 text_cleaning_tokens <- data %>%
26   tidytext::unnest_tokens(word, text)
27 text_cleaning_tokens$word <- gsub('[:digit:]]+', '', text_cleaning
  _tokens$word)
28 text_cleaning_tokens$word <- gsub('[:punct:]]+', '', text_cleaning
  _tokens$word)
29 text_cleaning_tokens <- text_cleaning_tokens %>% filter(!(nchar(
  word) == 1))%>%
30   anti_join(stop_words)
31 tokens <- text_cleaning_tokens %>% filter(!(word==""))
32
33
34 tokens <- tokens %>% mutate(ind = row_number())
35 tokens <- tokens %>% group_by(id) %>% mutate(ind = row_number())
  %>%
36   tidyr::spread(key = ind, value = word)
37 tokens[is.na(tokens)] <- ""
38 tokens <- tidyr::unite(tokens, text, -id, sep = "_" )
39 tokens$text <- trimws(tokens$text)
40
41 #create DTM
42 dtm <- CreateDtm(tokens$text,
43   doc_names = tokens$id,
44   ngram_window = c(1, 2))
45
46
47 # explore the basic frequency
48 tf <- TermDocFreq(dtm = dtm)

```

```

49 original_tf <- tf %>% select(term, term_freq, doc_freq)
50 rownames(original_tf) <- 1:nrow(original_tf)
51
52 # eliminate words appearing less than 2 times or in more than half
  of the doc
53 vocabulary <- tf$term[ tf$term_freq > 1 & tf$doc_freq < nrow(dtm) /
  2 ]
54
55 dtm = dtm
56
57 # run LDA
58 k_list <- seq(1, 20, by = 1)
59 model_dir <- paste0("models_", digest::digest(vocabulary, algo = "
  sha1"))
60 if (!dir.exists(model_dir)) dir.create(model_dir)
61 model_list <- TmParallelApply(X = k_list, FUN = function(k){
62   filename = file.path(model_dir, paste0(k, "_topics.rda"))
63
64   if (!file.exists(filename)) {
65     m <- FitLdaModel(dtm = dtm, k = k, iterations = 500)
66     m$k <- k
67     m$coherence <- CalcProbCoherence(phi = m$phi, dtm = dtm, M = 5)
68     save(m, file = filename)
69   } else {
70     load(filename)
71   }
72
73   m
74 }, export=c("dtm", "model_dir"))
75
76 # model tuning
77 # choosing the best model
78 coherence_mat <- data.frame(k = sapply(model_list, function(x) nrow
  (x$phi)),
79
80                                     coherence = sapply(model_list, function
  (x) mean(x$coherence)),
81                                     stringsAsFactors = FALSE)
82 ggplot(coherence_mat, aes(x = k, y = coherence)) +
83   geom_point() +
84   geom_line(group = 1)+
85   ggtitle("Best Topic by Coherence Score") + theme_minimal() +
86   scale_x_continuous(breaks = seq(1,20,1)) + ylab("Coherence")
87
88 # select models based on max average
89 model <- model_list[which.max(coherence_mat$coherence)][[ 1 ]]
90
91 # top 20 terms based on phi
92 model$top_terms <- GetTopTerms(phi = model$phi, M = 20)
93 top20_wide <- as.data.frame(model$top_terms)
94
95 # word, topic relationship

```

```

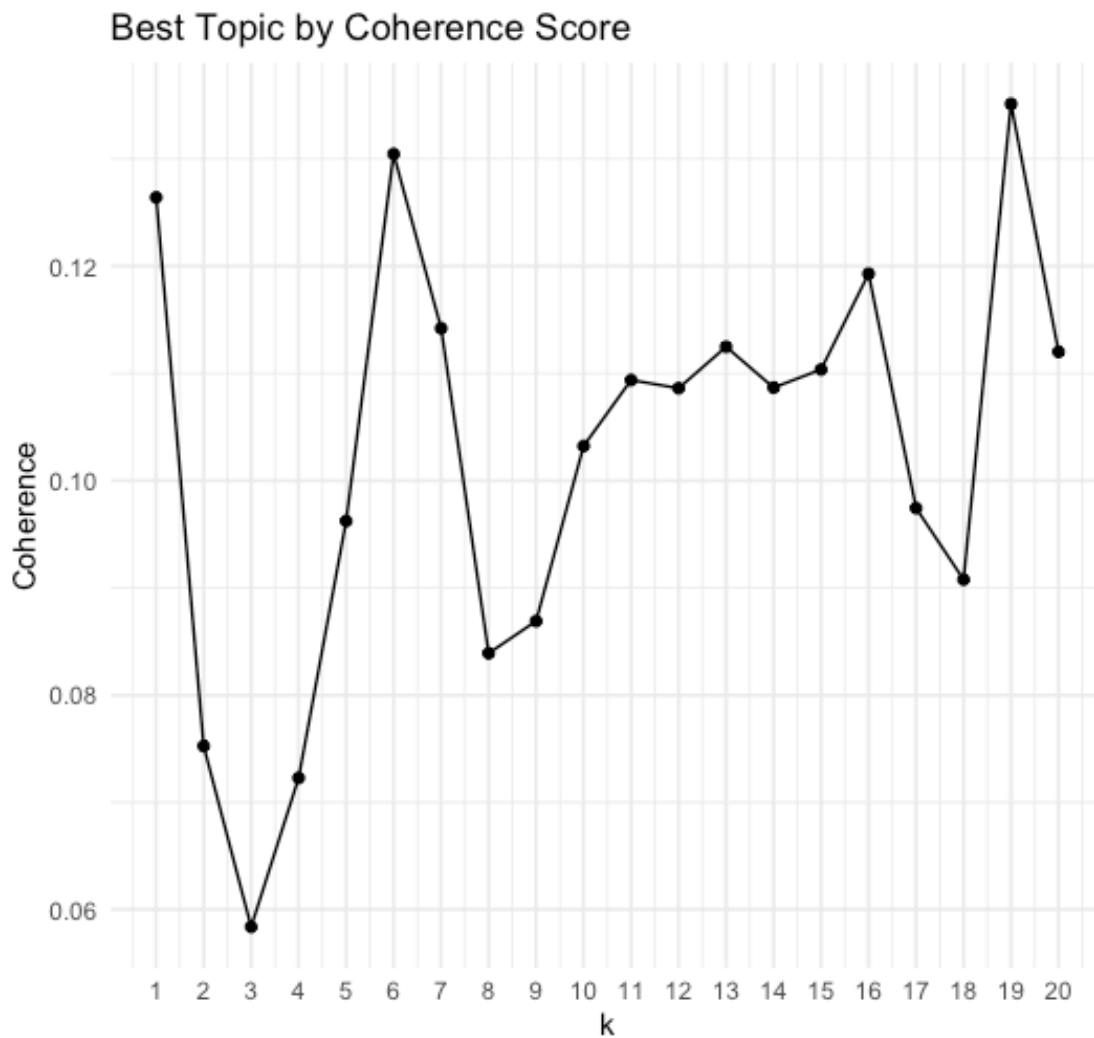
95 # looking at the terms allocated to the topic and their pr(word|
    topic)
96 allterms <- data.frame(t(model$phi))
97 allterms$word <- rownames(allterms)
98 rownames(allterms) <- 1:nrow(allterms)
99 allterms <- melt(allterms, idvars = "word")
100 allterms <- allterms %>% rename(topic = variable)
101 FINAL_allterms <- allterms %>% group_by(topic) %>% arrange(desc(
    value))
102
103
104 # topic, word, freq
105 final_summary_words <- data.frame(top_terms = t(model$top_terms))
106 final_summary_words$topic <- rownames(final_summary_words)
107 rownames(final_summary_words) <- 1:nrow(final_summary_words)
108 final_summary_words <- final_summary_words %>% melt(id.vars = c("
    topic"))
109 final_summary_words <- final_summary_words %>% rename(word = value)
    %>% select(-variable)
110 final_summary_words <- left_join(final_summary_words, allterms)
111 final_summary_words <- final_summary_words %>% group_by(topic, word)
    %>%
112   arrange(desc(value))
113 final_summary_words <- final_summary_words %>% group_by(topic, word
    ) %>% filter(row_number() == 1) %>%
114   ungroup() %>% tidyr::separate(topic, into = c("t", "topic")) %>%
    select(-t)
115 word_topic_freq <- left_join(final_summary_words, original_tf, by =
    c("word" = "term"))
116
117 # per-document-per-topic probabilities
118 theta_df <- data.frame(model$theta)
119 theta_df$document <- rownames(theta_df)
120 rownames(theta_df) <- 1:nrow(theta_df)
121 theta_df$document <- as.numeric(theta_df$document)
122 theta_df <- melt(theta_df, id.vars = "document")
123 theta_df <- theta_df %>% rename(topic = variable)
124 theta_df <- theta_df %>% tidyr::separate(topic, into = c("t", "topic"
    )) %>% select(-t)
125 FINAL_document_topic <- theta_df %>% group_by(document) %>%
126   arrange(desc(value)) %>% filter(row_number() == 1)
127
128 # visualising of topics in a dendrogram
129 model$topic_linguistic_dist <- CalcHellingerDist(model$phi)
130 model$hclust <- hclust(as.dist(model$topic_linguistic_dist), "ward.
    D")
131 model$hclust$labels <- paste(model$hclust$labels, model$labels[ ,
    1])
132 plot(model$hclust)
133
134

```

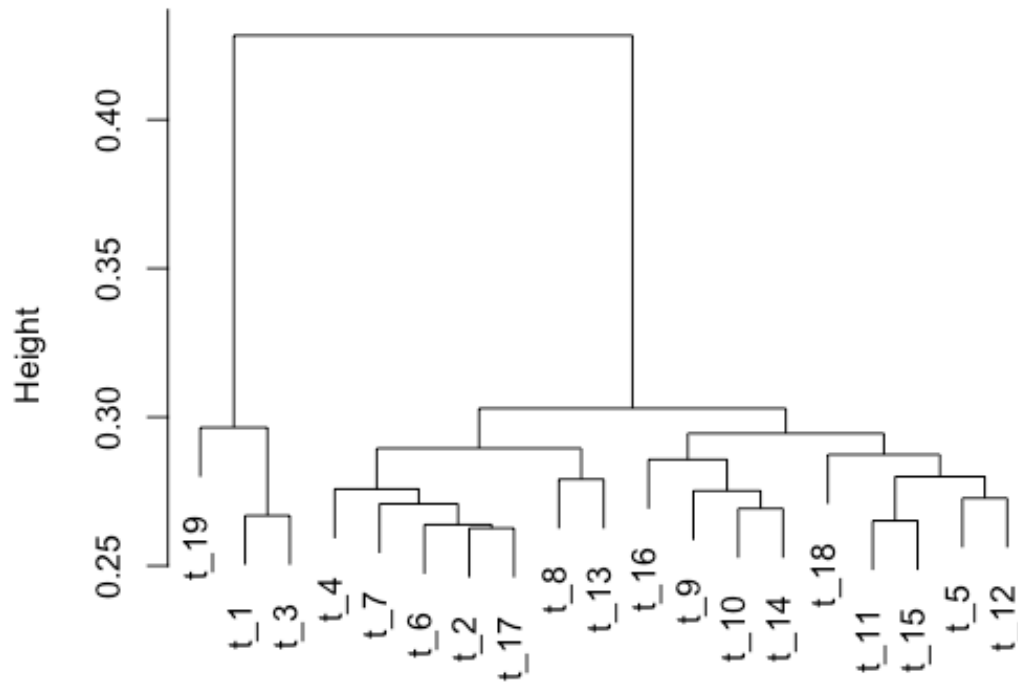
```

135 # visualising topics of words based on the max value
136 set.seed(1234)
137 pdf("result.pdf")
138 for(i in 1:length(unique(final_summary_words$topic)))
139 { wordcloud(words = subset(final_summary_words ,topic == i)$word,
140   freq = subset(final_summary_words ,topic == i)$value, min.freq
141   = 1,
142   max.words=200, random.order=FALSE, rot.per=0.35,
143   colors=brewer.pal(8, "Dark2"))}
144 dev.off()

```



Cluster Dendrogram



```
as.dist(model$topic_linguistic_dist)
hclust (*, "ward.D")
```

Figure 1: Word cloud of topic 3

