

DTS 301TC - Assessment 2 Project Report

Jiayuan Zhu
1931097

1 T1

Start by reading the xlsx file by column type.

```
1 col_types <- c("text", "text", "text", "text", "text", "text", "numeric", "numeric", "numeric", "numeric", "numeric")
2 data <- read_excel("TwitterDataset.xlsx", col_types = col_types)
```

Then, by checking the data we found three columns, **Favs**, **RTs**, and **Followers**, have missing values, set the missing values to 0.

```
1 data <- data %>%
2   mutate_at(c("Favs", "RTs", "Followers"), function(x)(
3     replace_na(x,0)))
```

1.1 T1-1

We created a new column **SumFavsRTs** by adding **Favs** and **RTs**, and then ranked the tweets according to the descending order of the values of a. The top ten rows are the top ten tweets. The Figure 1 shows the top ten tweets, with some columns hidden due to the small size of the position.

```
1 RankTweets <- data %>%
2   mutate(SumFavsRTs = Favs + RTs) %>%
3   arrange(desc(SumFavsRTs))
4
5 RankTweets[1:10, ]
```

[illegible]

Figure 1: Top Ten Tweets

1.2 T1-2

We started by selecting the three columns of **User Name**, **Nickname**, and **Followers**. We then removed the duplicate rows as each user may have posted more than one tweets, and ranked the users according to the reverse order of the values of **Followers**. The top ten rows are then taken to get the top ten users, as shown in Figure 2.

```
1 RankUsers <- data %>%
2   select("User Name", "Nickname", "Followers") %>%
3   filter(!duplicated(Nickname)) %>%
4   arrange(desc(Followers))
5
6 RankUsers[1:10, ]
```

```
# A tibble: 10 x 3
```

	User Name`	Nickname	Followers
	<chr>	<chr>	<dbl>
1	WEREVERTUMORRO	werevertumorro	7232508
2	Joel McHale	joelmchale	3804447
3	Questlove Gomez	questlove	3627104
4	Gabrielle Union	itsgabrielleu	3043936
5	Alfredo Flores	AlfredoFlores	2742081
6	MALUMA	maluma	2525716
7	Matthew Ziff	MatthewZiff	2139467
8	Fernanda Brum	PraFeBrum	2128396
9	Maxene Magalona	maxenemagalona	2031651
10	CristiandelaFuente	iamdelafuente	1965579

Figure 2: Top Ten Users

1.3 T1-3

Firstly, we get the Hour from the column **Date** and create a new column **Hour**, then we group all the rows according to the Hour and find the number of rows **Count** in each group. The Figure 3 shows the number of tweets posted in 24 hours.

```
1 TimeTweets <- data %>%
2   mutate(Hour = str_split(Hour, ":", simplify=T)[,1]) %>%
3   group_by(Hour) %>%
4   summarise(Count = n())
5
6 ggplot(TimeTweets) +
7   geom_line(aes(x = Hour, y = Count, group = 1), color='
8     #82B0D2') +
9   geom_point(x = TimeTweets$Hour, y = TimeTweets$Count,
10    color='#82B0D2') +
11   labs(x = 'Hour',
12    y = 'Number of Tweets',
13    title = " ") +
14   theme(plot.title = element_text(hjust = 0.5, size = 15))
15
16   ,
17   axis.text.x = element_text(size = 10),
18   axis.text.y = element_text(size = 10),
19   axis.title = element_text(size = 15))
```

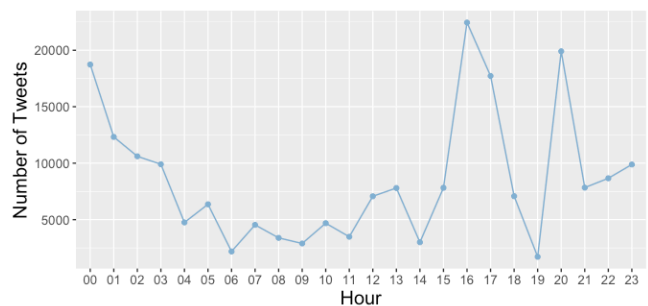


Figure 3: Number of Tweets Posted in a 24 hours

1.4 T1-4

I'd like to draw a figure to show the number of tweets posted from different states in the US. Firstly, we need to find the state to which the coordinates of each tweet belong. Here, I used package *sf* to do it. Get states map from it and change it to a spatial dataframe, then spatial join with our coordinates.

```
1 # Get the states map, turn into sf object
2 US <- st_as_sf(map("state", plot = FALSE, fill = TRUE))
3
4 # Make it a spatial dataframe, using the same coordinate
  system as the US spatial dataframe
5 CoordsData <- st_as_sf(data, coords = c("Longitude", "
  Latitude"), crs = st_crs(US))
6
7 # Perform a spatial join
8 sf_use_s2(FALSE)
9 StateTweets <- st_join(CoordsData, US)
```

Secondly, group each tweet according to the states they belong to and count the number of tweets in each group to get column **Count**.

```
1 StateTweets <- StateTweets %>%
2   drop_na(ID) %>%
3   group_by(ID) %>%
4   summarise(Count = n()) %>%
5   arrange(ID)
```

Finally, the number of tweets is presented on a map of the US states, from light to dark indicating the number of tweets from most to least, as shown in Figure 4.

```
1 us <- map_data("state")
2
3 ggplot() + geom_map(data=us, map=us,
4   aes(x=long, y=lat, map_id=region),
5   fill="ffffff", color="ffffff", size
6   =0.15) +
7   geom_map(data=StateTweets, map=us,
8     aes(fill=Count, map_id=ID),
9     color="ffffff", size=0.15) +
10   scale_fill_continuous(low='thistle2',
11     high='darkred',
12     guide='colorbar') +
13   labs(x=NULL, y=NULL) +
14   theme(axis.ticks = element_blank()) +
15   theme(axis.text = element_blank())
```

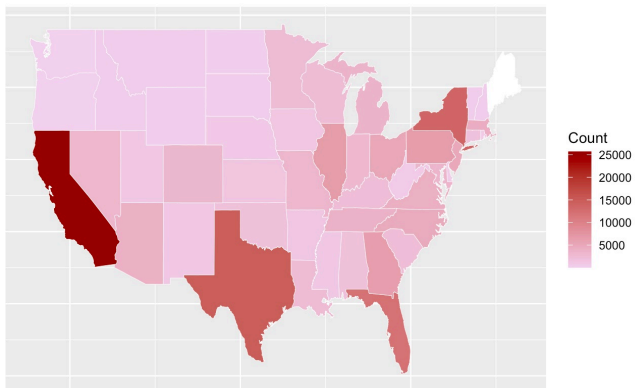


Figure 4: Number of Tweets Posted from Different States

2 T2

2.1 T2-1

To facilitate data mining, a function has been created for removing URLs, Emojis, mentions, hashtags and so on. Then it is applied on **Tweet content**. The Figure 5 shows some samples of tweet content

```
1 clean_tweets <- function(x) {
2   x %>%
3     # Remove URLs
4     str_remove_all("(?f|ht)(tp)(s?)(:|/)(.*)[.]/(.*)")
5     %>%
6     # Remove Emojis
7     str_remove_all('[:emoji:]') %>%
8     # Remove mentions
9     str_remove_all("@[[:alnum:]]{4,}") %>%
10    # Remove hashtags
11    str_remove_all("#[[:alnum:]]{4,}") %>%
12    # Replace "&" character reference with "and"
13    str_replace_all("&", "and") %>%
14    # Remove punctuation
15    str_remove_all("[[:punct:]]") %>%
16    # Remove "RT: "
17    str_remove_all("RT: ") %>%
18    # Replace any newline characters with a space
19    str_replace_all("\\\\n", " ") %>%
20    # Make everything lowercase
21    str_to_lower() %>%
22    # Remove multiple spaces by single space
23    str_squish()
24  }
25
26 d_clean <- data %>%
27   mutate_at(c("Tweet content"), clean_tweets)
28
29 head(d_clean$`Tweet content`)
```

[1] "wind mph nne barometer in rising slowly temperature "f rain today in humidity"
 [2] "pausa pro café antes de embarcar no próximo vôo trippolisontheroad danipolisviája pause for"
 [3] "good morning morning saturday diner vt breakfast nucorpsocadetsring ring college"
 [4] "recordstoredayus toms music trade"
 [5] "egg in a muffin rocket baby bakery in waumatosá wi"
 [6] "shouldve gave the neighbor a buzz iv got ice cream and moms baked goodies"

Figure 5: Samples of tweet content

2.2 T2-2

I'm going to do topic modeling in T3, so I'm going to remove the numbers, split into tokens and remove stop words for column **Tweet content** next. Some samples of tokens are shown in Figure 6.

```
1 d_token <- d_clean %>%
2   # Remove numbers
3   mutate(`Tweet content` = gsub("[[:digit:]]", "", `Tweet
4     content`)) %>%
5   # Split into tokens
6   unnest_tokens(word, `Tweet content`) %>%
7   # Remove stop words
8   anti_join(stop_words)
```

3 T3

3.1 Introduction

3.1.1 Background. Topic modelling is a technique designed to extract potential information from large datasets [1]. Topic models show enormous promise as a means of collecting potential insight


```

10 # Remove the duplicate
11 DTM = DTM[unique(DTM$), ]

```

Next, I used function **FindTopicsNumber** in package **ladatuning** to find the best number of topic in this task. The Figure 7 shows the performances of choosing different number K of topics between 4 to 10. It indicates that the performance gets better as K increases and still tends to increase. However, if $K > 10$, the computational effort increases and the visualisation is not favourable. So in the next step we set $K = 10$.

```

1 # Find the best K
2 result <- ldatuning::FindTopicsNumber(
3   DTM,
4   topics = seq(from = 4, to = 12, by = 2),
5   control = list(seed = 12345),
6   verbose = TRUE
7 )
8
9 FindTopicsNumber_plot(result)

```

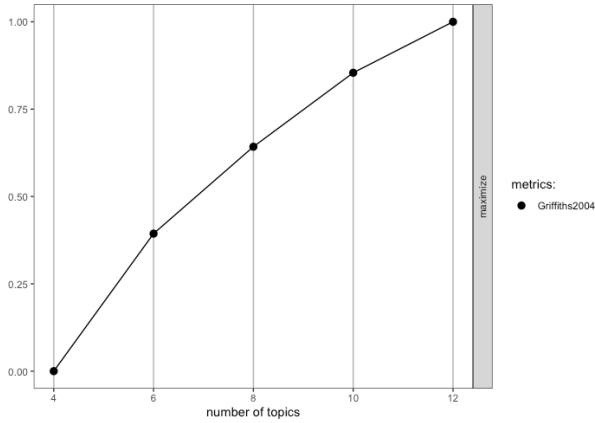


Figure 7: Top 15 Words of Each Topic

Next, the Gibbs Sampline and VEM methods of LDA were used to fit the data set respectively and get two models.

```

1 # Number of topics
2 K <- 10
3
4 set.seed(12345)
5 # Compute the LDA model, inference via 1000 iterations of
6   Gibbs sampling
7 GibbsModel <- LDA(DTM, K, method="Gibbs", control=list(
8   iter = 1000, verbose = 200))
9
10 # Compute the LDA model using VEM approach
11 VEMModel <- LDA(DTM, K, method="VEM")

```

3.4 Evaluation

To find the best model, we calculated the perplexity of each model.

Perplexity is a measure of how well a probability distribution or probability model predicts a sample. The lower the perplexity, the better the model clusters. The perplexity is calculated as follows

$$\text{perplexity}(D) = \exp\left(-\frac{\sum \log p(w)}{\sum_{d=1}^M N_d}\right)$$

	Gibbs Sampline	Variational EM
Perplexity	1816.72	2054.069

Table 1: Perplexity of Two Models

where $\sum_{d=1}^M N_d$ is the sum of all words in the test set, namely the total length of the test set. $p(w)$ refers to the probability of occurrence of each word in the test set:

$$p(w) = p(z|d) * p(w|z)$$

where $p(z|d)$ represents the probability of occurrence of each topic in a document and $p(w|z)$ represents the probability of occurrence of each word in the dictionary under a topic.

```

1 perplexity(VEMModel, DTM)
2 perplexity(GibbsModel, DTM)

```

The result shown in Table 1 indicates that **GibbsModel** has the lower perplexity, meaning that it performed better. Thus, we used **GibbsModel** for the next steps.

Then, we calculate the words of each topic. For clearer visualisation, we plot the 15 words with highest beta value of each topic in Figure 8.

```

1 # apply auto tidy using tidy and use beta as per-topic -
2   per-word probabilities
3 topic <- tidy(GibbsModel, matrix = "beta")
4
5 # choose 15 words with highest beta from each topic
6 top_terms <- topic %>%
7   group_by(topic) %>%
8   top_n(15, beta) %>%
9   ungroup() %>%
10  arrange(topic, -beta)
11
12 # plot the topic and words for easy interpretation
13 plot_topic <- top_terms %>%
14   mutate(term = reorder_within(term, beta, topic)) %>%
15   ggplot(aes(term, beta, fill = factor(topic))) +
16   geom_col(show.legend = FALSE) +
17   theme(axis.text.x = element_text(size = 7),
18         axis.text.y = element_text(size = 5),
19         axis.title = element_text(size = 15)) +
20   facet_wrap(~ topic, scales = "free") +
21   coord_flip() +
22   scale_x_reordered()
23 plot_topic

```

In order to see more clearly the connections between the different topics, we have generated a page with the results of the model clustering for analysis, as shown in Figure 9.

```

1 topicmodels2LDAvis <- function(x, ...){
2   post <- posterior(x)
3   if (ncol(post[["topics"]]) < 3) stop("The model must
4     contain > 2 topics")
5   mat <- x@wordassignments
6   createJSON(
7     phi = post[["terms"]],
8     theta = post[["topics"]],
9     vocab = colnames(post[["terms"]]),
10    doc.length = slam::row_sums(mat, na.rm = TRUE),
11    term.frequency = slam::col_sums(mat, na.rm = TRUE)
12  )
13 }

```

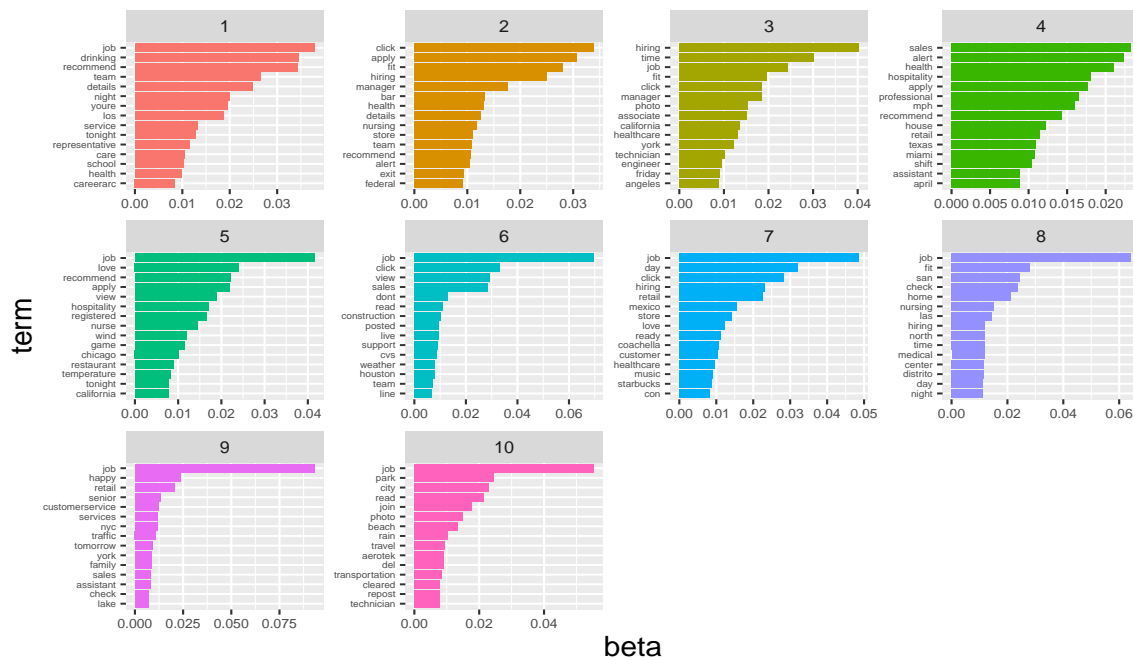


Figure 8: Top 15 Words of Each Topic

```

13 serVis(topicmodels2LDAvis(topicModel))
14

```

The visualisation results show that *job* is the most dominant topic in these topics amount the tweets, while the partitioning between different topics is not very obvious, indicating that the performance of the mod is bad and needs further improvement.

3.5 Conclusion

In this project, I used **LDA** for topic models for tweet dataset. However, the **LDA** has many shortcomings in topic modelling.

- The number of topics K is fixed that we need to tuning it for better performance. Also, as K increases, the consumption time increases.
- Topic distribution cannot capture correlations since it uses the Dirichlet distribution to model changes in topic proportions.
- The bag of words is used and the sentence structure information is not used.
- Most of the tweets are short texts, which have little semantic information and are sparse. **LDA** is not suitable it.

In order to improve the accuracy of topic modelling, structural information between sentences should also be used. Therefore, in the future there is a need to make more use of sentence structure information rather than just using the bag-of-words method.

REFERENCES

- [1] David M. Blei. 2012. Probabilistic Topic Models. *Commun. ACM* 55, 4 (apr 2012), 77–84. <https://doi.org/10.1145/2133806.2133826>
- [2] Yogesh Girdhar, Philippe Giguere, and Gregory Dudek. 2013. Autonomous adaptive underwater exploration using online topic modeling. In *Experimental Robotics*. Springer, 789–802.
- [3] Tom Griffiths. 2002. Gibbs sampling in the generative model of latent dirichlet allocation. (2002).
- [4] Liangjie Hong and Brian D Davison. 2010. Empirical study of topic modeling in twitter. In *Proceedings of the first workshop on social media analytics*. 80–88.
- [5] Suhyeon Kim, Haecheong Park, and Junghye Lee. 2020. Word2vec-based latent semantic analysis (W2V-LSA) for topic modeling: A study on blockchain technology trend analysis. *Expert Systems with Applications* 152 (2020), 113401.
- [6] Lin Liu, Lin Tang, Wen Dong, Shaowen Yao, and Wei Zhou. 2016. An overview of topic modeling and its current applications in bioinformatics. *SpringerPlus* 5, 1 (2016), 1–22.
- [7] Brendan O'Connor, Ramnath Balasubramanian, Bryan R Routledge, and Noah A Smith. 2010. From tweets to polls: Linking text sentiment to public opinion time series. In *Fourth international AAAI conference on weblogs and social media*.
- [8] Christos H Papadimitriou, Hisao Tamaki, Prabhakar Raghavan, and Santosh Vempala. 1998. Latent semantic indexing: A probabilistic analysis. In *Proceedings of the seventeenth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems*. 159–168.
- [9] Minghui Qiu, Feida Zhu, and Jing Jiang. 2013. It is not just what we say, but how we say them: Lda-based behavior-topic model. In *Proceedings of the 2013 SIAM international conference on data mining*. SIAM, 794–802.
- [10] Daniel Ramage, Evan Rosen, Jason Chuang, Christopher D Manning, and Daniel A McFarland. [n. d.]. Topic Modeling for the Social Sciences. ([n. d.]), 4.
- [11] Dionisios N Sotiropoulos, Chris D Kounavis, Panos Kourouthanassis, and George M Giaglis. 2014. What drives social sentiment? An entropic measure-based clustering approach towards identifying factors that influence social sentiment polarity. In *IISA 2014, The 5th International Conference on Information, Intelligence, Systems and Applications*. IEEE, 361–373.
- [12] Asbjørn Steinskog, Jonas Therkelsen, and Björn Gambäck. 2017. Twitter topic modeling by tweet aggregation. In *Proceedings of the 21st nordic conference on computational linguistics*. 77–86.
- [13] Jianshu Weng, Ee-Peng Lim, Jing Jiang, and Qi He. 2010. Twitterrank: finding topic-sensitive influential twitterers. In *Proceedings of the third ACM international conference on Web search and data mining*. 261–270.

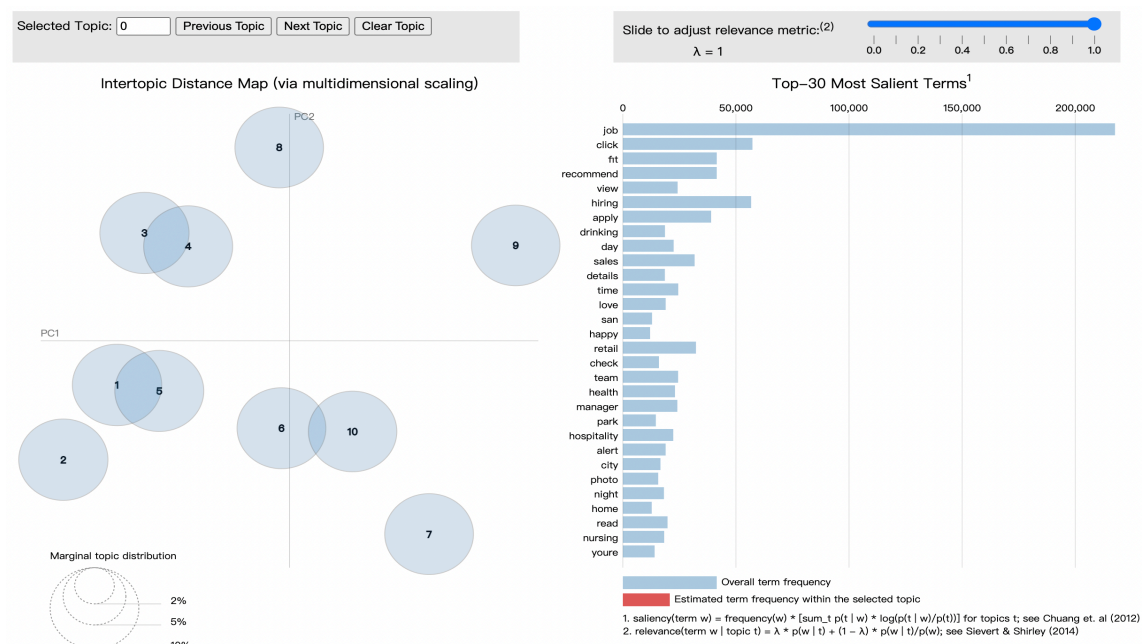


Figure 9: Topics