

Activity, Language, Use and Social Interactions in an Online Community

Dania Ismadi

8/05/2020

Contents

| | | |
|----------|---|-----------|
| 1 | Activity and language on the forum over time | 1 |
| 1.1 | Analysis of the activity of participants | 1 |
| 1.2 | Analysis of linguistic variables | 2 |
| 2 | Language used by different groups | 4 |
| 2.1 | Description of threads | 4 |
| 2.2 | Difference in language between groups | 4 |
| 2.3 | Language in threads over time | 5 |
| 3 | Social networks online | 7 |
| 4 | Conclusion | 9 |
| 5 | Appendix | 10 |
| 5.1 | Tidying Data | 10 |
| 5.2 | Additional Tables and Graphs | 11 |
| 5.2.1 | Part 1 | 11 |
| 5.2.2 | Part 2 | 12 |
| 5.3 | R Code | 18 |
| 5.3.1 | Part 1 | 18 |
| 5.3.2 | Part 2 | 24 |
| 5.3.3 | Part 3 | 28 |

1 Activity and language on the forum over time

This data is read from [webforum.csv](#) which is a file containing all the posts on an online forum. The thread ID, author ID, date, word count, and a number of language variables measured through proportions are used to describe each post. My personal sample of 20,000 posts contained quite a few anonymous authors. I decided to remove all these authors and their posts from the data for consistency as it would have been inaccurate to use these values for certain calculations. After doing so, I discovered that the data set spanned 10 years, from January 6, 2002 until December 31, 2011. During this time, there were a total of 2395 authors who posted in 600 different threads. The following subsections will go into a further exploration of the data set by analysing the activity of participants and the linguistic variables in the online forum. Note that all graphs or plots mentioned are included in the appendix.

1.1 Analysis of the activity of participants

First, we will examine the relationship between authors, the number of posts and word count. With 95% confidence, we can say that the average number of posts for each authors ranges from approximately 7 to 9 posts. However, the distribution is extremely right skewed as seen in the histogram. Most authors actually only post once and 75% of authors post less than 5 times. We see the same right-skewed distribution with word count. The average number of words for each authors ranges from approximately 701 to 941 words but 75% of authors wrote less than 515 words in total.

A t.test was performed in order to examine if authors who post more have a higher word count on average than authors who post less. Since the t.test requires an approximately normally distributed data and our data is extremely right-skewed, we have to modify the data a little bit. First, we will only include data among active users (those who post more than a certain amount) and exclude those who post a significant amount more than everybody else. A log transformation is then performed which results in an approximately normal distribution as shown in the histogram. We define authors who post the most as those who are in the top 10 percentile of this data and authors who post the least as those who are in the bottom 10 percentile of the this data. After performing a t.test, we find that there is no evidence that is strong enough to suggest that those who have post more have a higher word count on average than those who post the least. Thus, for the rest of the report, the most active authors are defined as those with the most posts and the least active authors as those with the least amount of posts. This is because from the t.test, we see that there is no reason to suggest that authors with high word counts interact with the forum more through posting than those that do not.

We can say with 95% confidence that the average number of threads authors participate in is between 4 and 5 threads but this is not representative of the population because the distribution is extremely right skewed. Most authors only participate in 1 thread. However, from high the t-statistic and low p-value obtained from performing a t.test, there is strong evidence to suggest that the most active authors participate in more threads on average than the authors who are least active.

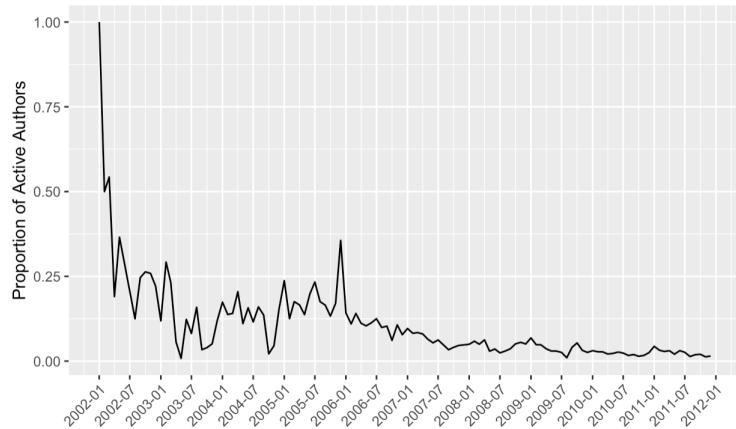


Figure 1: Time Series of the Proportion of Active Authors Over Time

Figure 1 is a time series of the proportion of active authors out of all the authors who ever interacted with the forum up until that time. We notice a generally decreasing trend, suggesting that most authors do not remain active on the thread for long. Most authors are usually only active for one day. With 95% confidence, we can say that authors are active for 209 to 247 days on average but the distribution is right skewed as 75% of authors are active for 197 days or less. Interesting enough, we cannot conclude the more active authors stay active in the forum for a longer time on average than less active authors.

We define an author's responsiveness as the average duration between their subsequent posts. We can say that with 95% confidence, authors take 81 to 101 days of break on average between subsequent postings. So since we could not conclude that more active authors stay active on the forum for long, then are they more responsive? This seems like a reasonable conclusion to make but from performing a t.test, the evidence is not strong enough to suggest that this is the case.

We know that most authors are not active on the forum for long but to get a high level view of the overall activity of participants over time, a time series graph of the number of posts posted per month onto the forum can be plotted (shown in figure 2). I have chosen to plot per month instead of per day as plotting posts per day leads to a messy graph and a trend would be difficult to observe. Furthermore, plotting per year is too low of a resolution as it would only yield us 10 data points. Therefore, plotting per month is the best option.

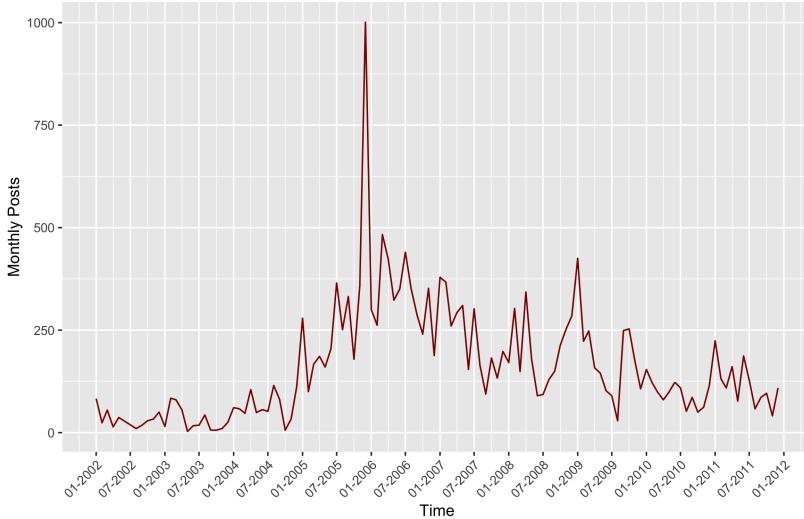


Figure 2: Number of Monthly Posts Over Time

From analysing the time series in figure 2, we can see that the number of posts start of low and gradually increase until it reaches a sharp spike in 2005. More specifically, it reaches an all time high of 1001 posts in October 2005. After, the number of posts begin to decrease. We also notice that 2005 is an outlier as the number of posts stay below 500 posts in other years. This trend is observed in the graph of the decomposition of this time series. Thus, it is clear that the activity of participants is not constant. However, excluding 2005, the change in the level activity is quite gradual.

1.2 Analysis of linguistic variables

The highest proportion of language variables present in posts belong to analytic, clout, authenticity and tone. The most significant out of those four variables is analytic, referring to words that express analytical thinking. This suggests that the forum might be geared more towards professional, academic or thought-provoking conversations rather than shallow, casual exchanges. A heat map of the proportion of language variables used over time in the forum is shown in figure 3. This heat map includes all threads in the web forum. Furthermore, I chose to plot the average proportion by year because plotting by month or day yields a very messy visualization as some months and days have no posts.

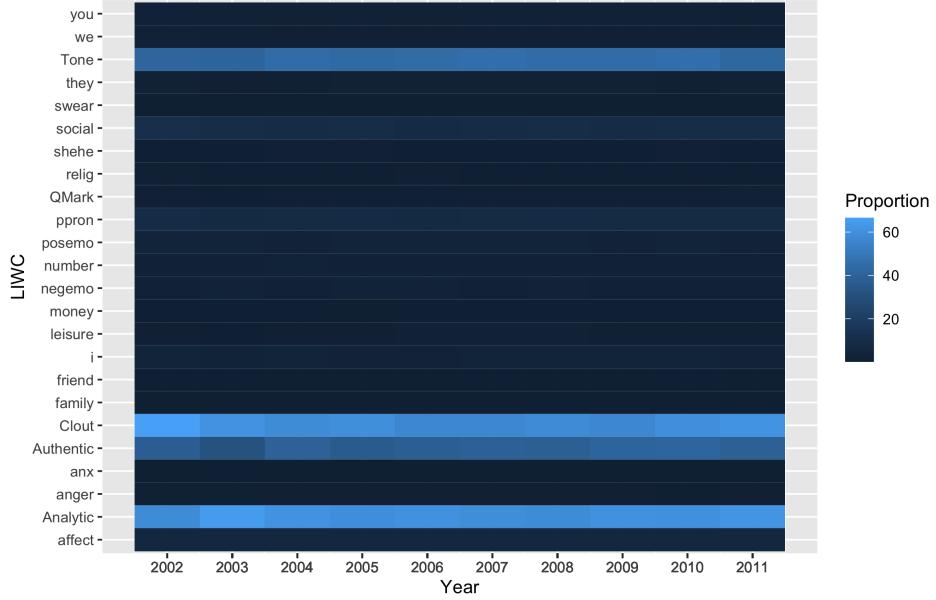


Figure 3: Heat Map of the Language Variables Proportions Over Time

From figure 3, we are able to clearly see the difference in proportions. Analytic, clout, authenticity, tone are significantly higher in proportion than the and other language variables. We also notice that the proportion of language variables generally stay consistent throughout the years. Again, this might suggest that the forum is quite niche regarding the types of conversations that is being discussed.

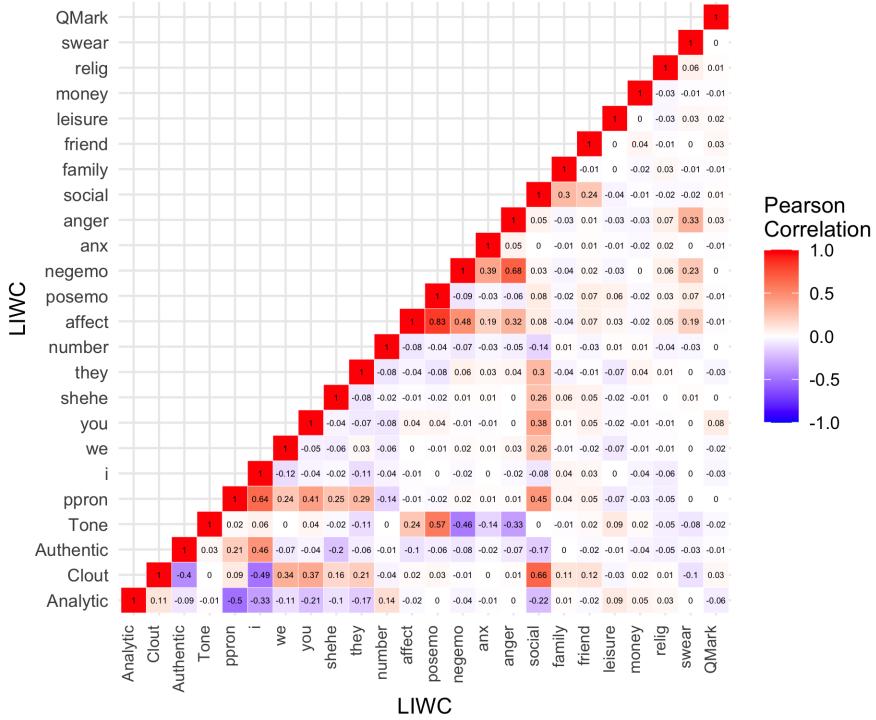


Figure 4: Correlation Matrix between Language Variables

Figure 4 shows a correlation matrix of the language variables in the form of a heatmap. We can observe that positive emotion and affect (words that express sentiment) are highly correlated with a value of 0.83. This suggests

that most posts that express an opinion or view in the forum are generally quite positive as well. Some other significant correlations to note are between anger and negative emotion (0.68) and between social and clout (0.66). Anger and negative emotion is quite self explanatory. As for social and clout, it appears that in this forum, authors express more power and force as words refer to social processes increase.

2 Language used by different groups

There are 600 different threads in this web forum. This section will go into a further exploration into the language used in these threads.

2.1 Description of threads

The distribution for the number of posts in each thread is right skewed as seen in the histogram. Most threads have 20 posts, and although the average number of posts is 32, 75% of threads have less than 34 posts. In order to satisfy the condition that the data has to be approximately normally distributed for a t.test, I excluded threads that are not active and threads that have significantly more posts than the rest. A log transformation is then performed resulting in approximately normally distributed data as seen in the histogram. More active threads are defined as the threads that belong in the top 10 percentile in this data and the less active threads are defined at the threads that belong in the bottom 10 percentile. From performing a t.test, there is strong evidence to suggest that threads that have more posts are higher in word count on average than threads with the less posts. Thus, we know that threads in this forum are not usually spammed with just one-line posts like how a digital messaging thread (i.e., Whatsapp, iMessage, Messenger) between people might look.

On average, there are 18 authors (17 to 19 authors with 95% confidence) per thread though most threads have 16 authors. The distribution is right skewed which means there are very few threads with large amounts of authors involved, however the maximum is 145 authors. There is strong evidence to conclude that more active threads have more contributing authors on average than less active threads.

Most threads are active for 3 days only, and 75% of threads are active for less than 112 days. However, the few threads that are active for longer are active for an extremely long time- the maximum is 2767 days. This pulls the average for the duration of active threads up to 178 days. There is strong evidence to suggest that more active threads are on average active for longer than less active threads. However, there is no evidence to suggest that more active threads have more frequent postings on average than less active threads. Thus, this suggests that for more active threads, the number of posts is less likely due to spam but due to actual involvement from the community. Most threads are updated with new posts every day and for 75% of posts the average break between posts is less than 10 days.

2.2 Difference in language between groups

To see if there is a difference in language variables used in different groups, I decided to compare between the most active and least active threads as perhaps, we might be able to observe some kind of pattern where certain threads that express a certain sentiment are more popular than others. A visualisation of this as a bar chart is shown in figure 5. At a glance, we can see that threads with more posts are more analytic, authentic, and revolve more around views and opinions (affect) while threads with the least posts express more force (clout), express more emotion (tone) and refer to social processes. However, to get a more accurate conclusion of the difference in language variables, t.tests were conducted for each variable. From the tests, we can attempt to make conclusions whether more popular threads express more of a certain sentiment on average than less popular threads. Some significant conclusions observed are that more popular threads express more negative emotion and anger on average than least popular threads. This suggests that perhaps it is more likely you will find posts that express strong negative opinion about a certain topic in more popular threads on average than less popular threads. Meanwhile, less popular threads express more force, power, emotion and words that refer to social processes (clout, tone and social). Referring to our previous section, we concluded that clout and social are highly correlated and now, we see that these variables appear more in less popular threads. Nevertheless, from the bar chart, it is clear that in both groups, the most popular language variables are consistent- analytic, clout, authentic and tone, as we also previously concluded.

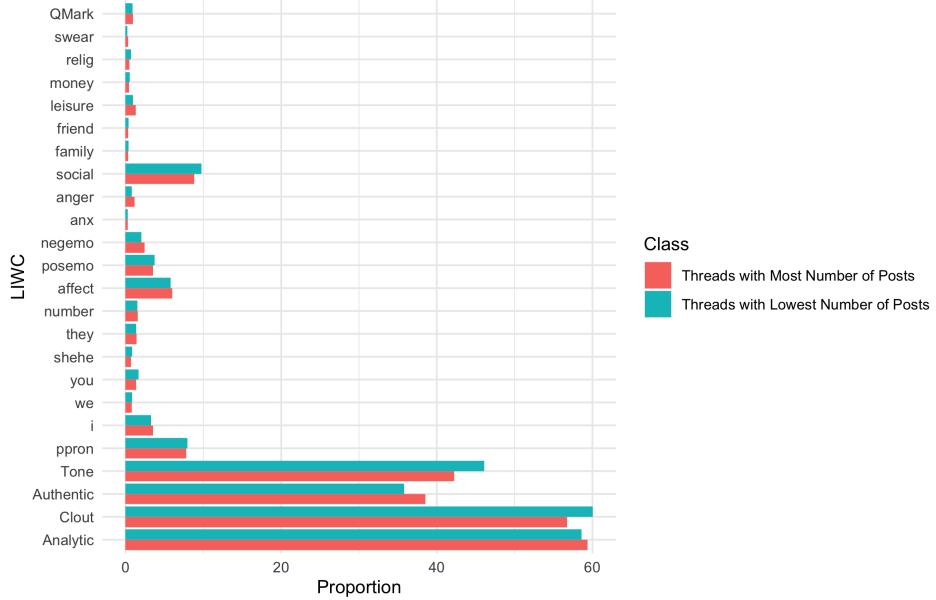


Figure 5: Bar Chart Comparing LIWC Between Most Active and Least Active Threads

2.3 Language in threads over time

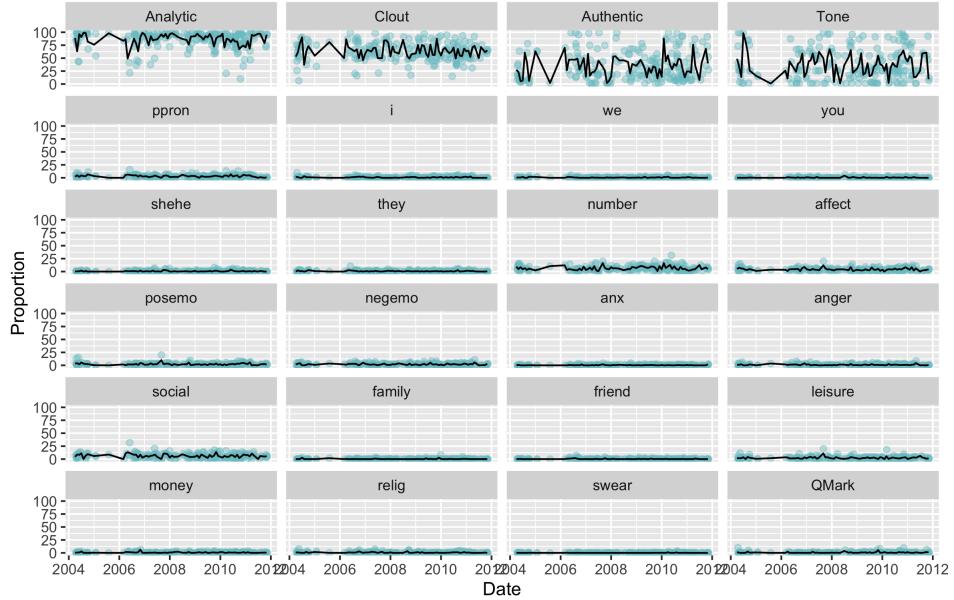


Figure 6: Comparing Language Variables for Thread 127715

To explore the question whether or not language in threads change over time, we will explore the top 6 threads with the most number of posts (Threads 252620, 127715, 472752, 145223, 283958, and 532649). I chose to do this because perhaps it increases the chance for a larger variety of different views, thoughts and opinions to be expressed so we might be able to see more observable changes than viewing a thread with less posts. To do this, I graphed a time series of the average proportion of each language variable used by month. I also mapped the average proportion of each language variable per day shown as the blue points in the graph. This makes it easy to observe how the data is distributed. The graph shown for the thread with thread ID 127715, the thread with the second most number of posts, is shown in figure 6.

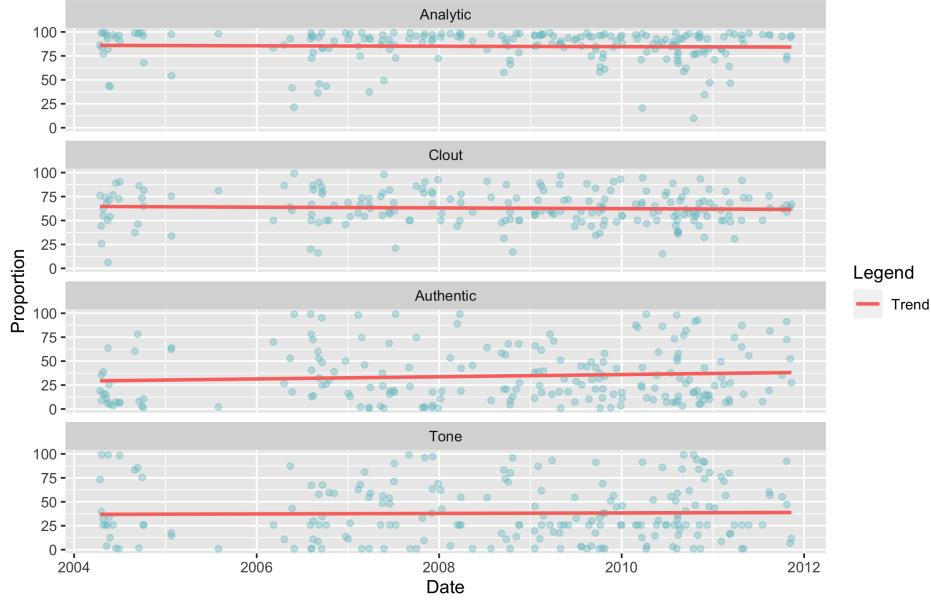


Figure 7: Comparing Language Variables for Thread 127115 with Trend Line

We see that in figure 6, most of the language variables actually stay constant (i.e., family, friend, leisure, QMark and may more) while variables like analytical thinking, clout, authenticity and tone are highly fluctuating showing that language for this thread does change from time to time. However, does this difference persist? The fluctuation makes it very difficult to see a trend over time so I plotted a trend line over the distribution to see if the difference does persist, focusing on the top 4 variables (analytic, clout, authentic and tone). This is shown in in figure 7. We see that the trend line is almost horizontal for all 4 variables. Thus, although these language variables fluctuate a lot, the difference does not persist over time. This suggests that for this thread, the sentiment expressed more or less remains the same during its life. From our sample, we see a similar trend for threads where authors post consistently over time- i.e., in threads 472752 and 532649. All graphs for each thread are included in the appendix in section 4.

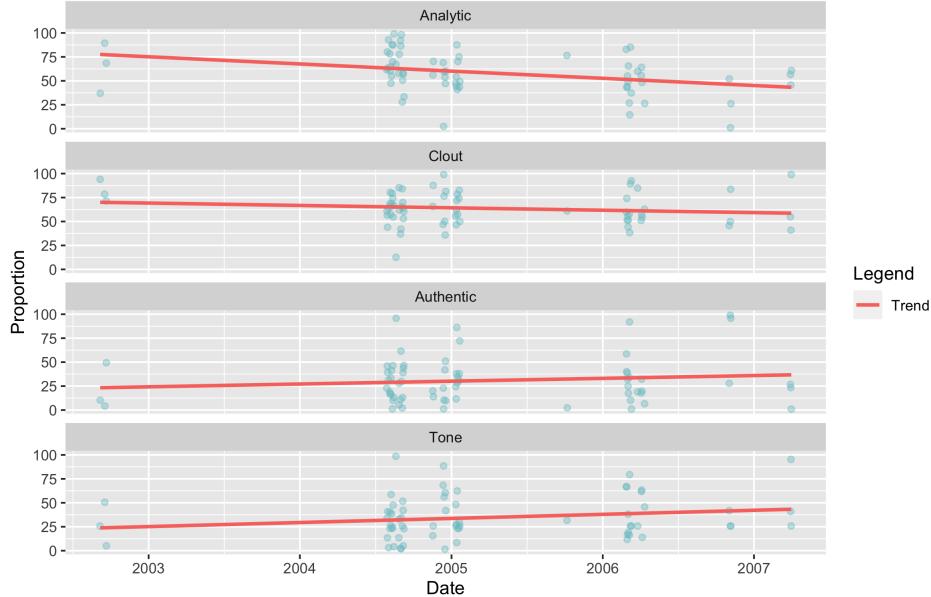


Figure 8: Comparing Language Variables for Thread 145223 with Trend Line

On the other hand, when we analyse threads that have large gap of time between postings, we see a slight difference. An example is shown for thread 145223 in figure 8. We see that for this particular thread, whilst analytical thinking and clout observed a decrease, authenticity and emotional tone saw an increase. Similar trends are observed with thread 252620 and 283958 where we see a slight increase or decrease in language variables. Thus, we can see that from this sample, it is more likely for a difference in language to persist in threads where there are less frequent postings. When threads observe more frequent postings (the duration between subsequent posts are short), the difference in language used in threads rarely persist over time though they do fluctuate. This suggests that for these threads more or less express the same sentiment for the duration that they are active.

3 Social networks online

Participants who communicate on the same network at the same time is a social network. For this analysis, I chose to examine the social network from March, April and May 2005 as this was when we started to see an increase in the number of people posting on the forum.



Figure 9: Social Network for March 2005

In the month of March, we see that there are 7 active threads resulting in 7 complete graphs. In these threads, 61 authors participated resulting in 440 connections. The largest clique size during this month is 16 which represents the thread with the most authors participating during March. During this month, all authors are connected to at least one another person so the result is a connected graph. This graph is shown in figure 9. More important authors are represented as bigger vertices while less important authors are smaller. From the histogram of the number of degrees against frequency, we see that the distribution is log normal which is typical of a social network.

There are lots of players with few connections and less players with more connections. The furthest connection between the two most distant people in this graph is 4 and the average path length is 2, thus this graph is very well connected. When observing the number of degrees, closeness centrality measure, betweenness centrality measure and eigencentrality, we notice that author 51994 has the top measurements. This author also belongs in the largest clique size. This suggests that this person is a strong local influencer and act as a bridge between different authors as well. In figure 9, this author is placed in the middle and their vertex is reasonably sized which emphasizes his important. Based on the high eigencentrality measurement, we can say that this person also has high quality connections as they are connected to other important authors.

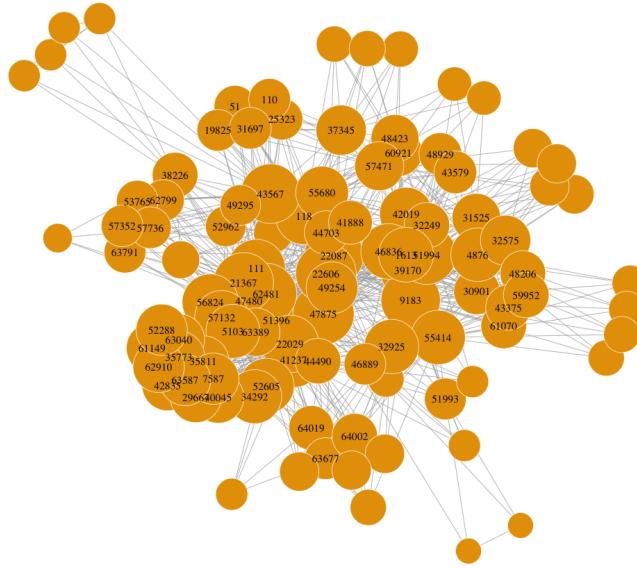


Figure 10: Social Network for March and April 2005

To see how this social network expands over time, we combine this graph with the graph in April. The resulting graph is shown in figure 10. We see that in this graph, the number of authors increase to 106 and the number of connections to 1011. The average path length stays the same but the furthest connection between the two most distant people decreased to 3 steps. This means that the authors is slightly more connected than before as perhaps some authors posted in many different threads with different authors. When observing the statistics, we see that author 47875 is high in the number of degrees and closeness centrality values. This means that this author has a lot of connections and is very well connected locally. Meanwhile author 435687 is high in betweenness centrality measure meaning that they act as a hub between other authors to connect. Finally, author 22606 has more quality connections due to its high eigencentrality value. However it is important to note that author 47875 follows a close second in these measures as well. This suggests that author 47875 is extremely well connected in this network. We see that this is because author 47875 is active in both March and April, and thus, was able to expand their own personal network. Whilst this person was not one of the most prominent figures in March, the statistics suggest that he was very important in April.

When extending our social network further to May, as expected, we see a large increase in the number of authors and connections to 133 and 1389 respectively. However, the average path length and diameter remains the same. Thus, it does not imply that the authors are significantly more well connected now than before. This social network is represented in figure 11. From analysing the statistics, we see that all authors that are high in the number of degrees, betweenness centrality, closeness centrality and eigencentrality measures have been active in all 3 months. This makes sense as there is a higher chance for them to interact with different authors if they are active throughout.

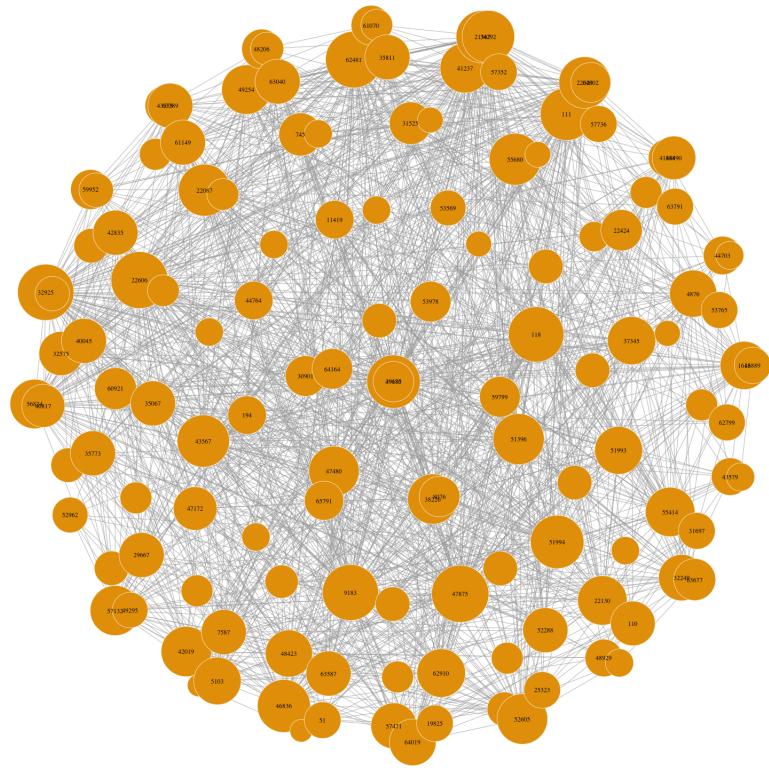


Figure 11: Social Network for March, April and May 2005

We also see that author 47875 remains an important player. However, we also notice that author 118 is also a top player during this time. Author 118 is the only author that appears as one of the top 6 players in each individual month and also when we built the network for March and April combined. This suggests that author 118 has consistently been a top player throughout March and April.

4 Conclusion

To summarise, activity from participants in this forum changes throughout. There are some periods where the forum experience increased activity and times where it does not but overall, the forum seems to have reached its peak in 2005 and has been decreasing in activity ever since. However, language, when measured by language variables, stays quite consistent in the forum throughout the years with the greatest proportions hailing from words that express analytical thinking, power and impact (clout), authenticity and emotional tone. There are some slight differences in the language used between different groups of threads but we generally see the same trend where language remains consistent through the duration of activity on the thread, suggesting that threads do not stray far from the sentiment they were originally expressing. Finally, from analysing the period from March 2005 to May 2005, we see typical attributes of a social network in this forum such as the fact that there are few people with very many connections while the majority have fewer connections. We were also able to point out a few key players during this time.

5 Appendix

5.1 Tidying Data

Importing the libraries:

```
1 library(readr)
2 library(ggplot2)
3 library(dplyr)
4 library(tidyverse)
5 library("tidyrr")
6 library(reshape2)
```

Before running each code block, here is how I set up my data:

```
1 # clear working environment
2 rm(list=ls())
3
4 # create individual data
5 set.seed(30990971)
6 webforum <- read.csv("webforum.csv")
7 # create a sample of 20,000 rows
8 webforum <- webforum[sample(nrow(webforum), 20000), ]
9 # remove all anonymous authors.
10 webforum = webforum[webforum$AuthorID != -1,]
11 # change Date column to date data type
12 webforum$date = as.Date(webforum$date, "%Y-%m-%d")
```

5.2 Additional Tables and Graphs

5.2.1 Part 1

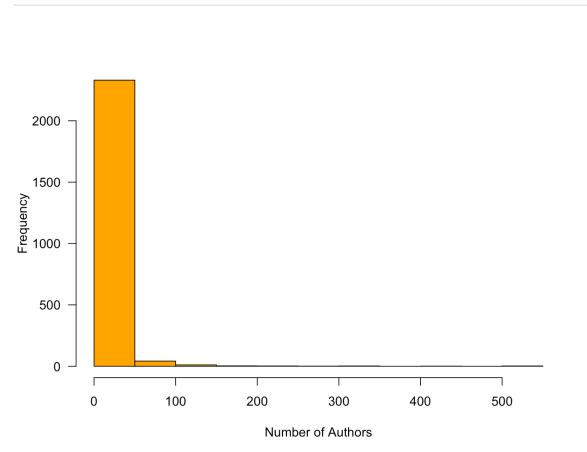


Figure 1: Histogram of the Number of Authors and the Total Number of Posts

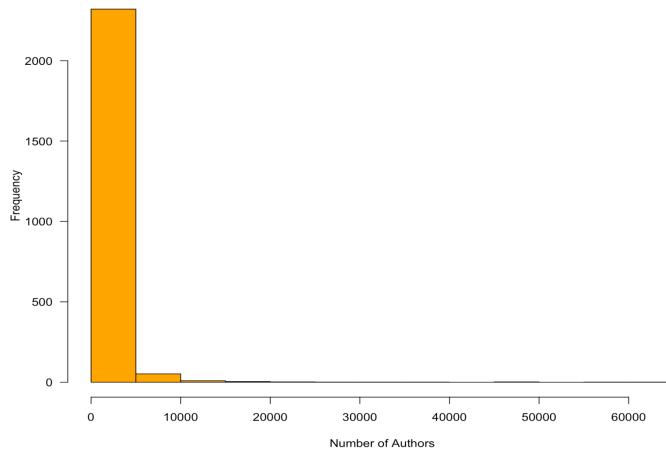


Figure 2: Histogram of the Number of Authors and the Total Word Count

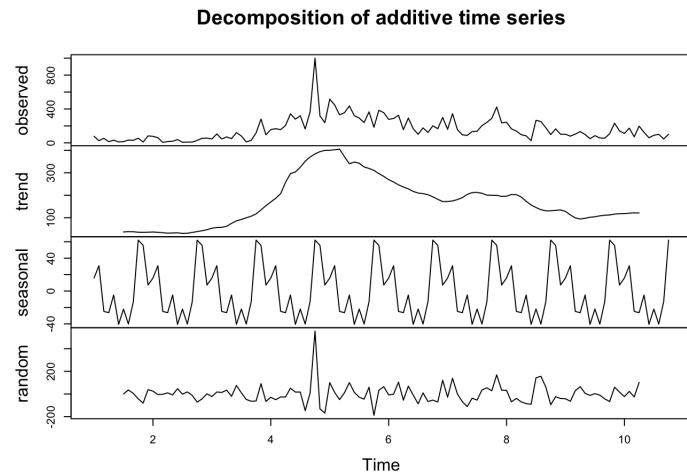


Figure 3: Decomposition of the Number of Monthly Posts Over Time

5.2.2 Part 2

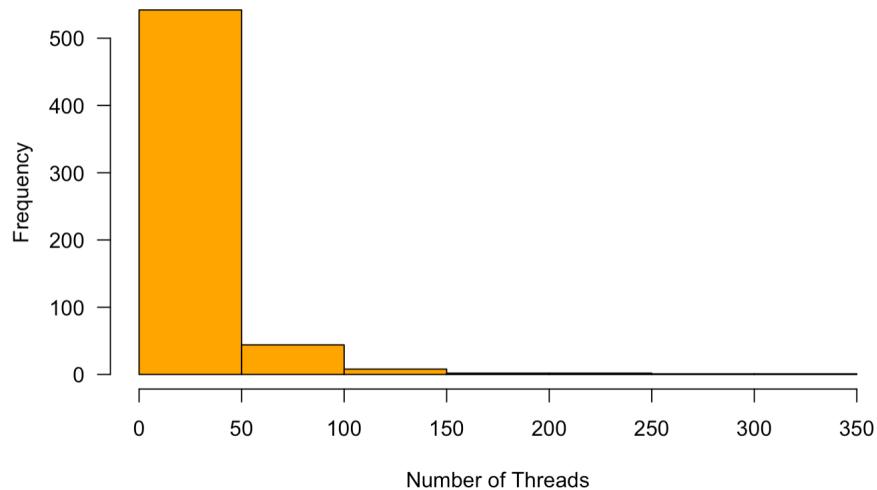


Figure 4: Histogram of the Number of Threads and Number of Posts

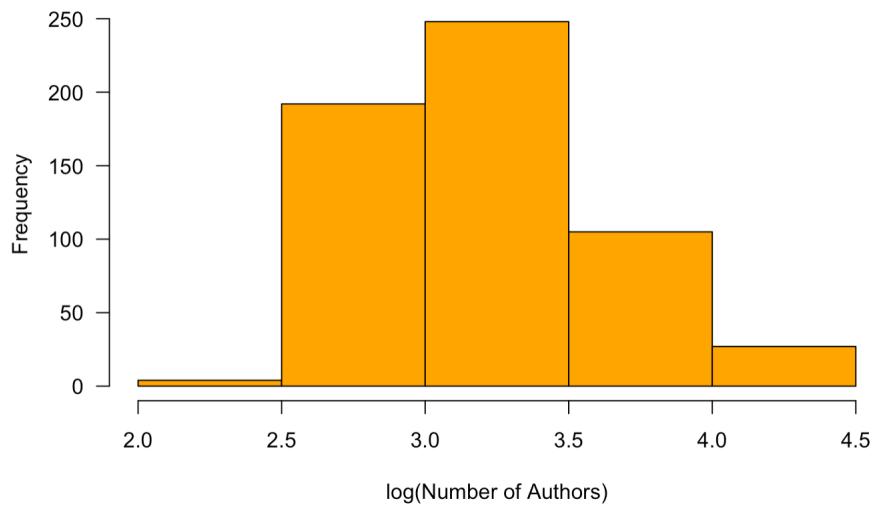


Figure 5: Histogram of the $\log(\text{Number of Threads})$ and $\log(\text{Number of Posts})$

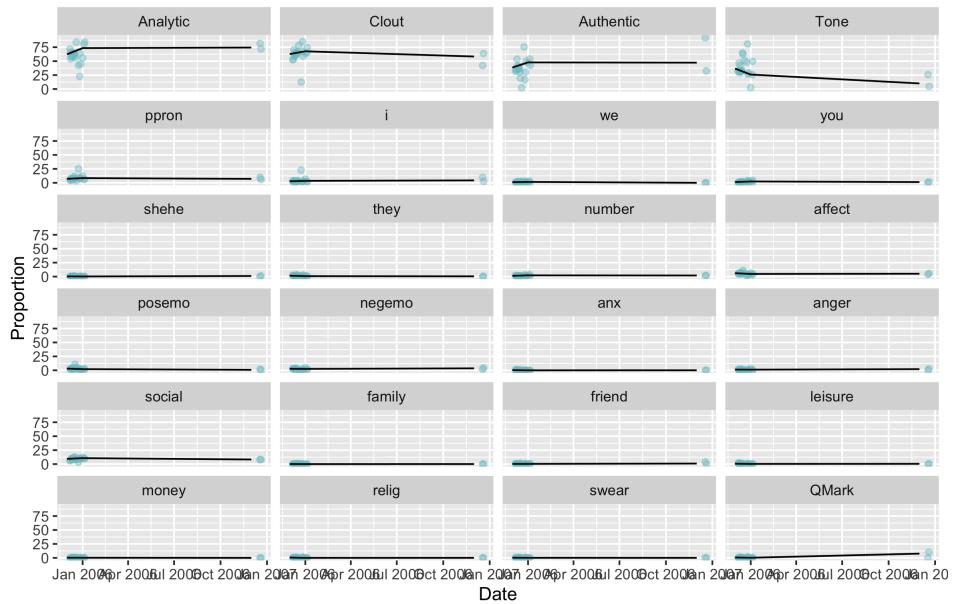


Figure 6: Comparing Language Variables for Thread 252620

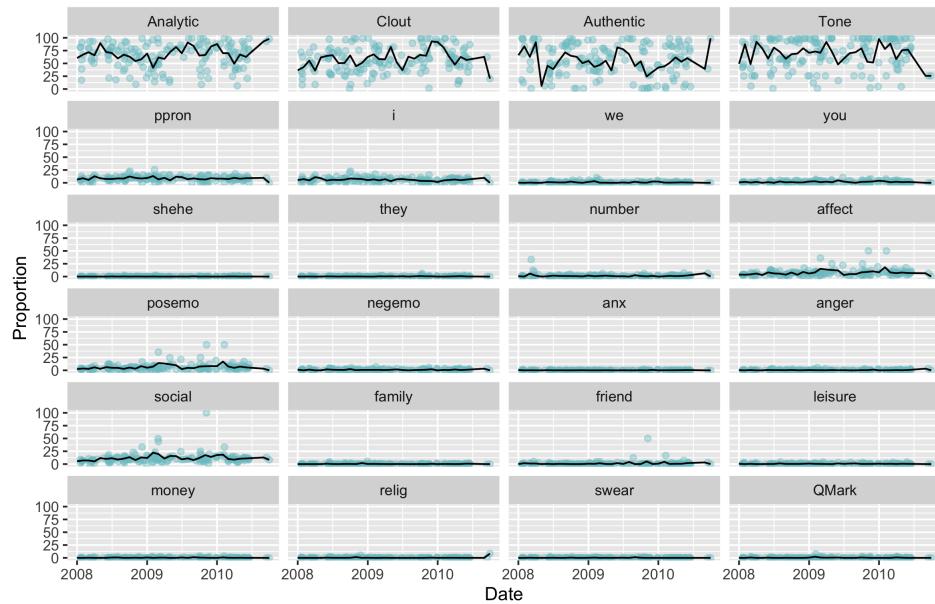


Figure 7: Comparing Language Variables for Thread 472752

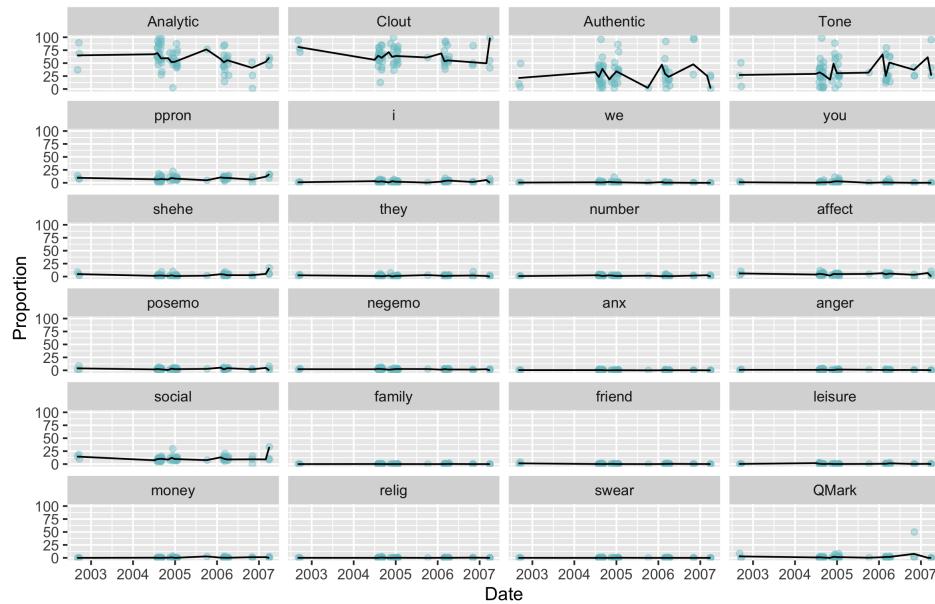


Figure 8: Comparing Language Variables for Thread 145223

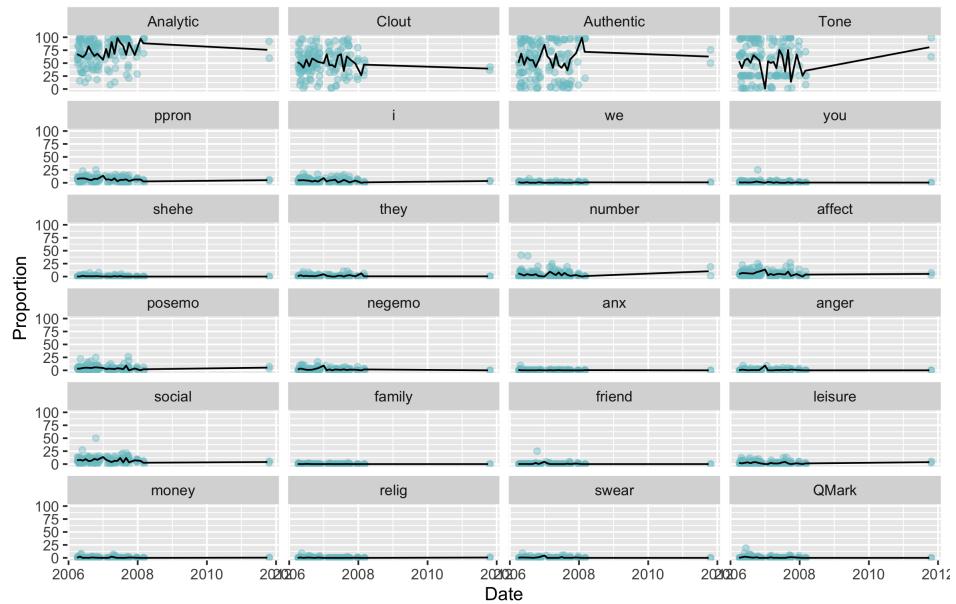


Figure 9: Comparing Language Variables for Thread 283958

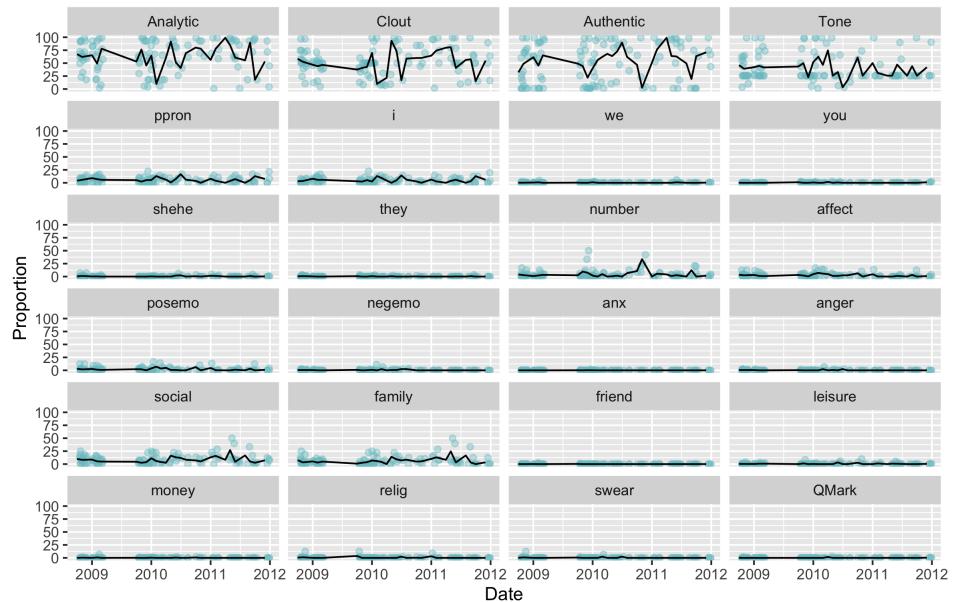


Figure 10: Comparing Language Variables for Thread 532649

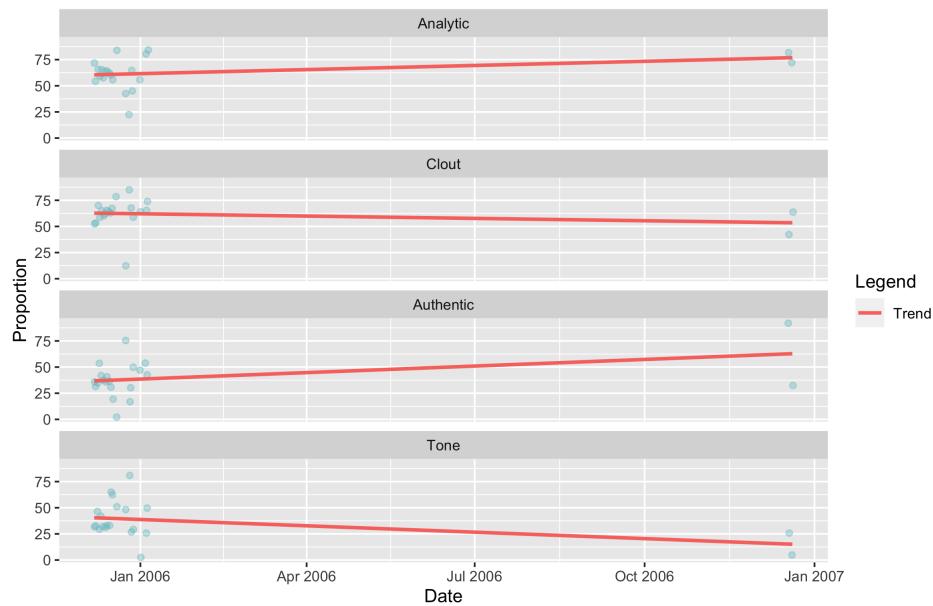


Figure 11: Comparing Language Variables for Thread 252620 with Trend Line

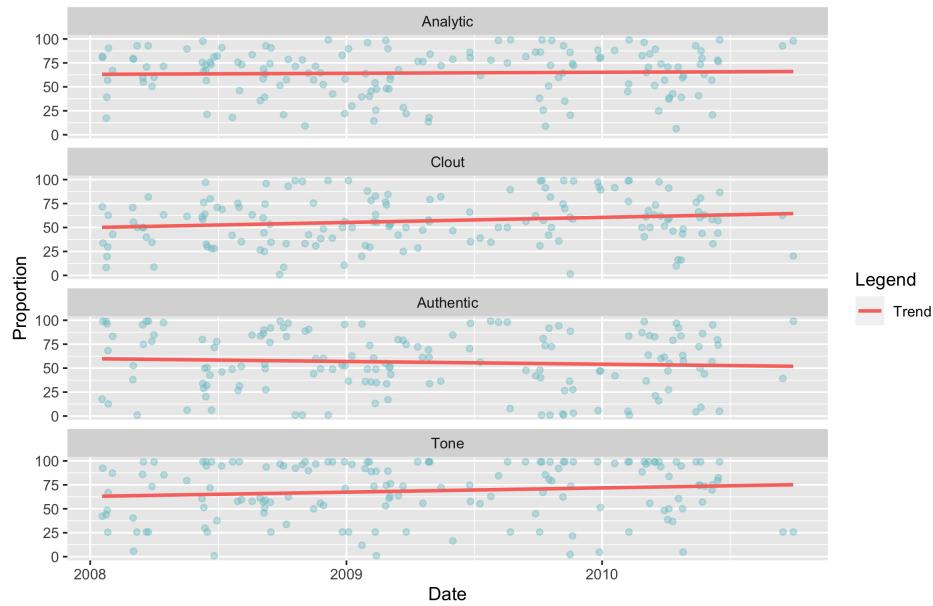


Figure 12: Comparing Language Variables for Thread 472752 with Trend Line

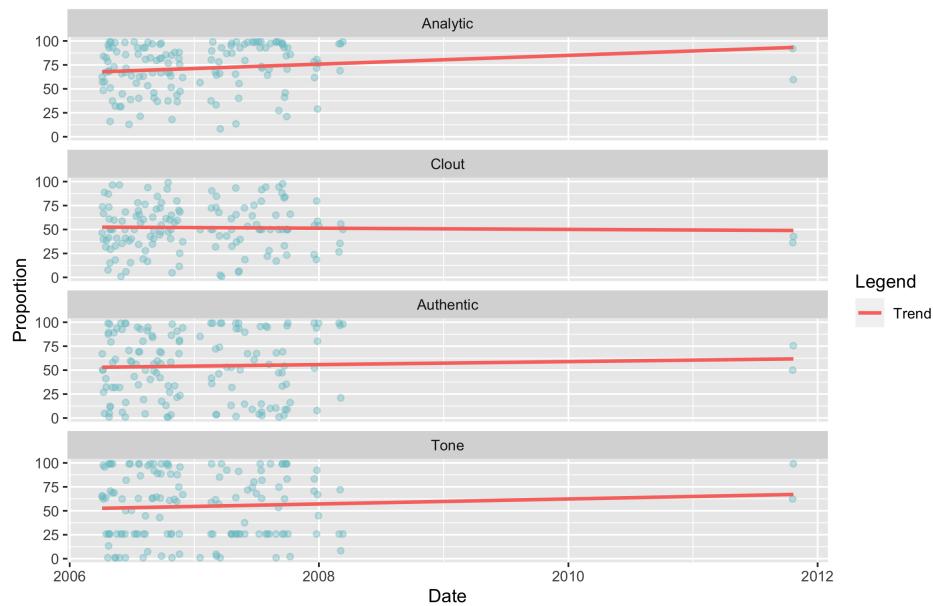


Figure 13: Comparing Language Variables for Thread 283958 with Trend Line

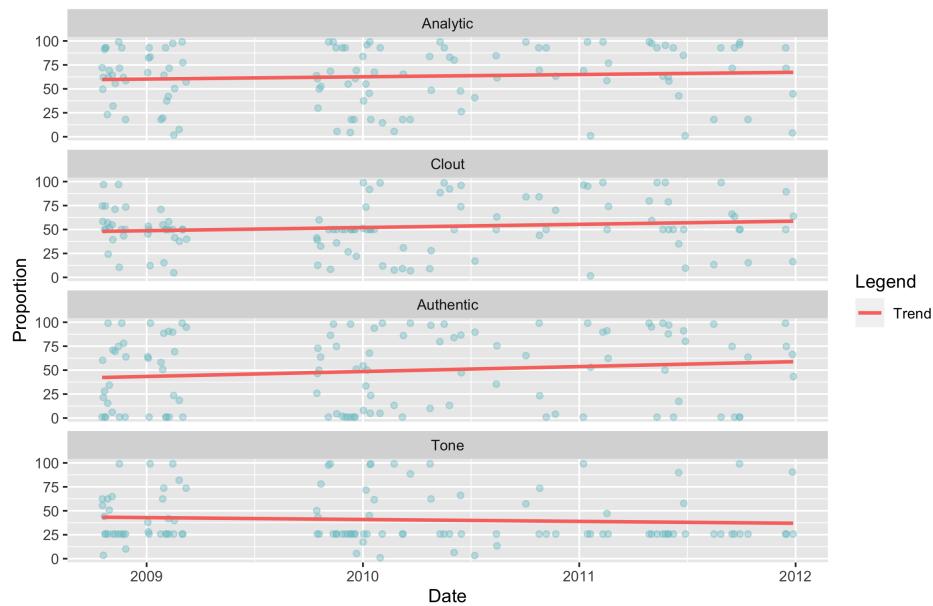


Figure 14: Comparing Language Variables for Thread 532649 with Trend Line

5.3 R Code

5.3.1 Part 1

Finding out the date range of the entire data:

```
1 # Find the minimum date.
2 min.date = webforum[which.min(as.Date(webforum$Date, "%Y-%m-%d")),]
3
4 # Find the maximum date.
5 max.date = webforum[which.max(as.Date(webforum$Date, "%Y-%m-%d")),]
6
7 min.date
8 max.date
9 ''''
```

Counting the number of threads:

```
1 # Calculate the number of threads and how often they appear in the data set.
2 num.threads = by(webforum, list(webforum$ThreadID), function(df) nrow(df))
3
4 # Combine output from by function into dataframe.
5 num.threads = as.data.frame(as.table(num.threads))
6 # Rename columns
7 colnames(num.threads) = c("ThreadID", "Count")
8
9 # Organise num.threads from least count to largest count.
10 num.threads = num.threads[order(num.threads$Count),]
11
12 dim(num.threads)
13 head(num.threads)
14 tail(num.threads)
```

Finding out how many authors are involved.

```
1 # Calculate the number of authors and how often they appear in the dataset.
2 num.authors <- by(webforum, list(webforum$AuthorID), function(df) nrow(df))
3
4 # Combine output from by function into dataframe.
5 num.authors <- as.data.frame(as.table(num.authors))
6 # Rename columns
7 colnames(num.authors) <- c("AuthorID", "Count")
8
9 # Orgnaise num.authors from least count to largest count.
10 num.authors <- num.authors[order(num.authors$Count),]
11
12 head(num.authors, 10)
13 tail(num.authors, 30)
14 dim(num.authors)
```

Statistics for the relationship between authors and the number of posts:

```
1 # Using data from before. Find statistics and the distribution of the number of posts
  per author.
2 summary(num.authors$Count)
3 getmode(num.authors$Count)
4 hist(num.authors$Count, main="", xlab="Number of Authors", border="black", col="orange",
      las=1, breaks=10)
5 # confidence interval
6 t.test(num.authors$Count)
```

Statistics for the relationship between authors and the total word count:

```
1 wc.authors = by(webforum$WC, list(webforum$AuthorID), sum)
```

```

2 wc.authors = as.data.frame(as.table(wc.authors))
3 colnames(wc.authors) = c("AuthorID", "TotalWC")
4 wc.authors = wc.authors[order(wc.authors$TotalWC),]
5
6
7 head(wc.authors)
8 tail(wc.authors)
9 summary(wc.authors$TotalWC)
10 getmode(wc.authors$TotalWC)
11 hist(wc.authors$TotalWC, main="", xlab="Number of Authors", border="black", col="orange"
12   , las=1, breaks=10)
12 # confidence interval
13 t.test(wc.authors$TotalWC)

```

Fixing the right-skew distribution:

```

1 num.authors
2 dim(num.authors[num.authors$Count > 80,])
3 fix.skew = num.authors[num.authors$Count >= 60,]
4 fix.skew = fix.skew[fix.skew$Count < 110,]
5 hist(fix.skew$Count, main="", xlab="log(Number of Authors)", border="black", col="orange"
6   , las=1, breaks=10)
6 hist(log(fix.skew$Count), main="", xlab="log(Number of Authors)", border="black", col="orange"
7   , las=1, breaks=5)

```

Finding the most active and least active authors, and performing t.tests against word count:

```

1 x = as.numeric(quantile(log(fix.skew$Count), prob=0.90))
2 y = as.numeric(quantile(log(fix.skew$Count), prob=0.10))
3 fix.skew$Count = log(fix.skew$Count)
4 fix.skew
5
6 # Define authors who post more as those who post in the top 10%.
7 most.posts = fix.skew[fix.skew$Count >= quantile(fix.skew$Count, prob=0.90),]
8 # Define authors who post less as those who post bottom 10%.
9 least.posts = fix.skew[fix.skew$Count <= quantile(fix.skew$Count, prob=0.10),]
10
11 most.posts.wc = wc.authors[wc.authors$AuthorID %in% most.posts$AuthorID,]
12 least.posts.wc = wc.authors[wc.authors$AuthorID %in% least.posts$AuthorID,]
13
14 # Perform t.test to conclude significance.
15 t.test(most.posts.wc$TotalWC, least.posts.wc$TotalWC, alternative="greater")
16
17 more.active = most.posts
18 less.active = least.posts

```

Statistics for the number of threads authors participate in:

```

1 webforum.temp = webforum
2
3 # Group by AuthorID
4 num.threads = as.data.frame(as.table(by(webforum.temp$WC, list(webforum.temp$ThreadID,
5   webforum.temp$AuthorID), FUN=length)))
5 num.threads = na.omit(num.threads)
6 colnames(num.threads) = c("ThreadID", "AuthorID", "NumPosts")
7
8 sum.threads = num.threads %>% group_by(AuthorID) %>%
9   arrange(ThreadID) %>%
10  summarize(NumThreads = as.numeric(length(ThreadID)))
11
12 summary(sum.threads$NumThreads)
13 getmode(sum.threads$NumThreads)
14 # histograms

```

```

15 hist(sum$threads$NumThreads, main="", xlab="Days", border="black", col="orange", las=1,
16   breaks=10)
17
18 # confidence intervals
19 t.test(sum$threads$NumThreads)

```

T.test comparing average thread count between more active and less active authors.

```

1 less.active.threads = sum$threads[(sum$threads$AuthorID %in% less.active$AuthorID),]
2 more.active.threads = sum$threads[(sum$threads$AuthorID %in% more.active$AuthorID),]
3
4 t.test(more.active.threads$NumThreads, less.active.threads$NumThreads, alternative="greater")

```

Statistics for the proportion of authors that remain active.

```

1 # Analyse activity by month.
2 webforum.temp = webforum
3 webforum.temp$date = as.Date(webforum.temp$date, "%Y-%m-%d")
4 webforum.temp$MonthYear = as.Date(format(webforum.temp$date, "%Y-%m-01"))
5
6 # Group AuthorID by month to find total number of Authors that post per month.
7 new.authors = as.data.frame(as.table(by(webforum.temp$WC, list(webforum.temp$MonthYear,
8   webforum.temp$AuthorID), FUN=sum)))
8 new.authors = na.omit(new.authors)
9 new.authors = new.authors[,-3]
10
11 # Rename Column Names.
12 colnames(new.authors) = c("Date", "AuthorID")
13
14 # Find total number of authors that post per month.
15 total.authors = as.data.frame(as.table(by(new.authors$AuthorID, list(new.authors>Date),
16   FUN=length)))
16 colnames(total.authors) = c("MonthYear", "TotalAuthors")
17 # Change MonthYear Column to Date data type.
18 total.authors$MonthYear = as.Date(total.authors$MonthYear, "%Y-%m-%d")
19
20 # Plot total number of authors that post per month.
21 g = ggplot(total.authors, aes(MonthYear, TotalAuthors)) + geom_line() + xlab("") + ylab(
22   "Number of Authors")
23 g = g + scale_x_date(date_labels = "%Y-%m", date_breaks = "6 month")
24 g = g + theme(axis.text.x = element_text(angle = 45, hjust = 1))
25 g
26
27 # remove duplicated rows and only preserve the first row which is the first time Author
28 # appeared in dataset to find only new authors who post per month.
29 new.authors = new.authors %>% distinct(AuthorID, .keep_all=TRUE)
30
31
32 # sort data frame by date
33 new.authors = new.authors[order(new.authors$Date),]
34
35 # find number of new authors permonth.
36 new.authors = as.data.frame(as.table(by(new.authors$AuthorID, list(new.authors>Date),
37   FUN=length)))
38 # rename column names
39 new.authors[is.na(new.authors)] = 0
40 colnames(new.authors) = c("MonthYear", "NewAuthors")
41 # change MonthYear column to date data type
42 new.authors$MonthYear = as.Date(new.authors$MonthYear, "%Y-%m-%d")
43
44 # plot total number of new authors that post per month.
45 t = ggplot(new.authors, aes(MonthYear, NewAuthors)) + geom_line() + xlab("") + ylab("New
46   Authors")

```

```

42 t = t + scale_x_date(date_labels = "%Y-%m", date_breaks = "6 month")
43 t = t + theme(axis.text.x = element_text(angle = 45, hjust = 1))
44 t
45
46 # find the proportion of the total number of authors that post per month out of the
47 # total number of authors who have ever posted in the forum until then
48 total.authors.percentages = total.authors
49 proportions = character()
50 for (i in 1:dim(total.authors)[1]) {
51   proportions = c(proportions, total.authors.percentages[i, 2] / sum(new.authors[1:i, 2]
52     )))
53 }
54 # column bind
55 total.authors.percentages = cbind(total.authors.percentages, proportions)
56 colnames(total.authors.percentages) = c("MonthYear", "TotalAuthors", "Proportion")
57
58 # change to numeric
59 total.authors.percentages$Proportion = as.numeric(as.character(total.authors.percentages
60   $Proportion))
61 total.authors.percentages
62
63 # plot the proportion
64 p = ggplot(total.authors.percentages, aes(MonthYear, Proportion)) + geom_line() + xlab(""
65   " MonthYear") + ylab("Proportion of Active Authors")
66 p = p + scale_x_date(date_labels = "%Y-%m", date_breaks = "6 month")
67 p = p + theme(axis.text.x = element_text(angle = 45, hjust = 1))
68 p

```

Statistics for the duration of authors are active for on the forum and t.test comparing more active users against less active users.

```

1 webforum.temp = webforum
2 # convert to date
3 webforum.temp$date = as.Date(webforum.temp$date, "%Y-%m-%d")
4
5 # group by date and author
6 dur = as.data.frame(as.table(by(webforum.temp$WC, list(webforum.temp$date, webforum.temp
7   $AuthorID), FUN=length)))
8 dur = na.omit(dur)
9 dur = dur[,-3]
10 colnames(dur) = c("Date", "AuthorID")
11 # convert to date data type
12 dur$date = as.Date(dur$date, "%Y-%m-%d")
13
14 # convert to character
15 dur$AuthorID = as.character(dur$AuthorID)
16
17 # Find the date ranges, inclusive.
18 date.ranges = dur %>% group_by(AuthorID) %>%
19   arrange(Date) %>%
20   summarize(Duration = as.numeric(diff(range(Date)))+1)
21
22 date.ranges$AuthorID = as.numeric(date.ranges$AuthorID)
23 date.ranges = date.ranges[order(date.ranges$AuthorID),]
24 summary(date.ranges$Duration)
25 getmode(date.ranges$Duration)
26
27 hist(date.ranges$Duration, main="",
28       xlab="Days",
29       border="black",
30       col="orange",
31       las=1,
32       )

```

```

31     breaks=20)
32
33 # Calculate confidence intervals for these means.
34 t.test(date.ranges$Duration)
35
36 # Can we conclude that more active authors are active for a longer period of time on avg
37 # . on the forum than less active authors?
38 less.active.daterange = date.ranges[(date.ranges$AuthorID %in% less.active$AuthorID),]
39 more.active.daterange = date.ranges[(date.ranges$AuthorID %in% more.active$AuthorID),]
40 t.test(more.active.daterange$Duration, less.active.daterange$Duration, alternative="greater")

```

Statistics for the responsiveness of authors and t.test:

```

1 webforum.temp = webforum
2 # convert to date
3 webforum.temp$date = as.Date(webforum.temp$date, "%Y-%m-%d")
4
5 # group by date and author
6 dur = as.data.frame(as.table(by(webforum.temp$WC, list(webforum.temp$date, webforum.temp
7   $AuthorID), FUN=length)))
7 dur = na.omit(dur)
8 dur = dur[,-3]
9 colnames(dur) = c("Date", "AuthorID")
10 # convert to date data type
11 dur$date = as.Date(dur$date, "%Y-%m-%d")
12
13 # convert to character
14 dur$AuthorID = as.numeric(as.character(dur$AuthorID))
15
16 # find frequency in dur
17 dur.freq = as.data.frame(as.table(by(dur$date, list(dur$AuthorID), FUN=length)))
18 colnames(dur.freq) = c("AuthorID", "NumPosts")
19 dur.freq$AuthorID = as.numeric(as.character(dur.freq$AuthorID))
20 dur.freq = dur.freq[order(dur.freq$AuthorID),]
21
22 # Remove all authors that have posted only once as we cannot measure their
23 # responsiveness.
24 dur.freq = dur.freq[dur.freq$NumPosts!=1,]
25
26 # Remove all authors that have only posted once in dur.
27 dur = dur[(dur$AuthorID %in% dur.freq$AuthorID),]
28
29 # Find average duration between posts for each author.
30 avg.dur = dur %>% group_by(AuthorID) %>%
31   arrange(Date) %>%
32   summarize(avg = as.numeric(mean(diff(Date))))
33
34 summary(avg.dur$avg)
35 getmode(avg.dur$avg)
36 hist(avg.dur$avg, main="", xlab="Days", border="black", col="orange", las=1, breaks=10)
37
38 # Calculate confidence intervals.
39 t.test(avg.dur$avg)
40
41 less.active.avgdur = avg.dur[(avg.dur$AuthorID %in% less.active$AuthorID),]
42 more.active.avgdur = avg.dur[(avg.dur$AuthorID %in% more.active$AuthorID),]
43 t.test(more.active.avgdur$avg, less.active.avgdur$avg, alternative="less")

```

Plotting the number of monthly posts over time and decomposing the time series:

```

1 # Consider month and year only
2 webforum$MonthYear = as.Date(format(webforum$Date, "%Y-%m-01"))
3
4 posts.per.month = as.data.frame(as.table(by(webforum$WC, list(webforum$MonthYear), FUN=
  length)))
5 colnames(posts.per.month) = c("Date", "NumPosts")
6 # convert date column to date data type
7 posts.per.month$Date = as.Date(posts.per.month$Date, "%Y-%m-%d")
8
9 # find the variance
10 var(posts.per.month$NumPosts)
11
12 # plot
13 t = ggplot(posts.per.month, aes(Date, NumPosts)) + geom_line(color="red4") + xlab("Time")
  ) + ylab("Monthly Posts")
14 t = t + scale_x_date(date_labels = "%m-%Y", date_breaks = "6 month")
15 t = t + theme(axis.text.x = element_text(angle = 45, hjust = 1))
16 t = t
17 t
18
19 # save the time series ggplot graph
20 ggsave("~/Desktop/FIT3152_resources/posts_over_time.png", t, width = 20, height = 13,
  units = "cm")
21
22 posts.ts = ts(posts.per.month$NumPosts, frequency = 12)
23 posts.ts
24 # decomposition of time series
25 decomp = decompose(posts.ts)
26 plot(decomp)

```

Figure 3: Heat Map:

```

1 # Descriptive statistics
2 summary(webforum)
3
4 # plot heat map of linguistic variables vs. year where shades represent average
5 webforum$Year = as.numeric(format(webforum$Date, "%Y"))
6 liwc.year.avg = aggregate(webforum[6:29], list(webforum$Year), FUN=mean)
7 colnames(liwc.year.avg)[1] = "Year"
8 liwc.year.long.avg = gather(data = liwc.year.avg, key = LIWC, value = Proportion, -1)
9
10 # heat map
11 a = ggplot(data=liwc.year.long.avg, aes(x=Year, y=LIWC))
12 a = a + geom_tile(aes(fill=Proportion))
13 a = a + scale_x_continuous(labels=liwc.year.long.avg$Year, breaks=liwc.year.long.avg$Year)
14 a

```

Figure 4: Correlation matrix as heat map:

```

1 webforum$Date = as.Date(webforum$Date, "%Y-%m-%d")
2
3 # plot heat map of linguistic variables vs. year where shades represent sum
4 corr = round(cor(webforum[6:29]), digits=2)
5
6 # Get lower triangle of the correlation matrix
7 get_lower_tri<-function(cormat){
  cormat[upper.tri(cormat)] <- NA
  return(cormat)
}
# Get upper triangle of the correlation matrix
get_upper_tri <- function(cormat){

```

```

13     cormat[lower.tri(cormat)] <- NA
14     return(cormat)
15   }
16
17 # reorder matrix due to correlation coefficient
18 reorder_cormat <- function(cormat){
19 # Use correlation between variables as distance
20 dd <- as.dist((1-cormat)/2)
21 hc <- hclust(dd)
22 cormat <- cormat[hc$order, hc$order]
23 }
24
25 cormat = reorder_cormat(corr)
26 upper.tri = get_upper_tri(corr)
27 # melted
28 melted.cormat <- melt(upper.tri, na.rm = TRUE)
29 # Create a ggheatmap
30 g <- ggplot(melted.cormat, aes(Var2, Var1, fill = value)) +
31   geom_tile(color = "white") +
32   scale_fill_gradient2(low = "blue", high = "red", mid = "white",
33                         midpoint = 0, limit = c(-1,1), space = "Lab",
34                         name="Pearson\nCorrelation") +
35   theme_minimal() + # minimal theme
36   theme(axis.text.x = element_text(angle = 90, vjust = 0,
37                                     size = 8, hjust = 1)) +
38   coord_fixed()
39 g = g + geom_text(aes(label = value), color="black", size=1.5)
40 g = g + xlab("LIWC") + ylab("LIWC")
41 g

```

5.3.2 Part 2

Describing summary statistics for the number of posts per thread.

```

1 # Count the number of threads.
2 threads = as.data.frame(as.table(by(webforum$WC, list(webforum$ThreadID), FUN=length)))
3 colnames(threads) = c("ThreadID", "NumPosts")
4
5 # Order from least posts to the most posts.
6 threads = threads[order(threads$NumPosts),]
7 summary(threads$NumPosts)
8 boxplot(threads$NumPosts)
9 hist(threads$NumPosts, xlab="Number of Threads", border="black", col="orange", las=1,
       breaks=10)
10 getmode(threads$NumPosts)
11
12 # Fix the right skew.
13 fix.skew = threads[threads$NumPosts > 0,]
14 fix.skew = fix.skew[fix.skew$NumPosts < 80,]
15
16 fix.skew
17
18 hist(fix.skew$NumPosts, main="", xlab="log(Number of Authors)", border="black", col="orange",
       las=1, breaks=7)
19 hist(log(fix.skew$NumPosts), main="", xlab="log(Number of Authors)", border="black", col="orange",
       las=1, breaks=7)
20
21 # Define threads who are more active as those who are in the top 10%.
22 top.threads = fix.skew[fix.skew$NumPosts >= quantile(fix.skew$NumPosts, prob=0.90),]
23 # Define threads who are less active as those are in the bottom 10%.

```

```

24 bottom.threads = fix.skew[fix.skew$NumPosts <= quantile(fix.skew$NumPosts, prob=0.10),]
25 top.threads

```

Threads by word count:

```

1 # Count the number of threads.
2 wc.threads = as.data.frame(as.table(by(webforum$WC, list(webforum$ThreadID), FUN=sum)))
3 colnames(wc.threads) = c("ThreadID", "TotalWC")
4
5 # Order
6 wc.threads = wc.threads[order(wc.threads$TotalWC),]
7
8 summary(wc.threads)
9
10 # Can we conclude that threads with a larger number of posts have a higher word count
   than threads with a lower number of posts?
11 top.posts.wc = wc.threads[wc.threads$ThreadID %in% top.threads$ThreadID,]
12 bottom.posts.wc = wc.threads[wc.threads$ThreadID %in% bottom.threads$ThreadID,]
13
14 t.test(top.posts.wc$TotalWC, bottom.posts.wc$TotalWC, alternative="greater")

```

Summary statistics for the number of authors per thread and perform t-test for the most popular and least popular threads.

```

1 # Group by number of threads and authors.
2 author.threads = as.data.frame(as.table(by(webforum$WC, list(webforum$ThreadID, webforum
   $AuthorID), FUN=length)))
3 colnames(author.threads) = c("ThreadID", "AuthorID", "NumPosts")
4 author.threads = na.omit(author.threads)
5 author.threads = author.threads[,-3]
6 author.threads = as.data.frame(as.table(by(author.threads$AuthorID, list(author.threads$ThreadID),
   FUN=length)))
7 colnames(author.threads) = c("ThreadID", "NumAuthors")
8 author.threads = author.threads[order(author.threads$NumAuthors),]
9
10 summary(author.threads$NumAuthors)
11 getmode(author.threads$NumAuthors)
12 boxplot(author.threads$NumAuthors)
13 hist(author.threads$NumAuthors)
14 t.test(author.threads$NumAuthors)
15
16 # Can we conclude that more popular threads (by post) have more authors posting in it
   than less popular threads? Yes, quite significant.
17 top.posts = author.threads[author.threads$ThreadID %in% top.threads$ThreadID,]
18 bottom.posts = author.threads[author.threads$ThreadID %in% bottom.threads$ThreadID,]
19 t.test(top.posts$NumAuthors, bottom.posts$NumAuthors)

```

Describe the summary statistics for how long threads are usually active for. Perform t-test for the most popular and least popular threads.

```

1 activity.period = as.data.frame(as.table(by(webforum$WC, list(webforum$Date, webforum$ThreadID),
   FUN=length)))
2 activity.period = na.omit(activity.period)
3 activity.period = activity.period[,-3]
4 colnames(activity.period) = c("Date", "ThreadID")
5
6 # Convert to date data type.
7 activity.period$Date = as.Date(activity.period$Date, "%Y-%m-%d")
8 activity.period$ThreadID = as.character(activity.period$ThreadID)
9

```

```

10 # Calculate date ranges.
11 date.ranges = activity.period %>% group_by(ThreadID) %>%
12   arrange(Date) %>%
13   summarize(Duration = as.numeric(diff(range(Date)))))
14
15 date.ranges$ThreadID = as.numeric(date.ranges$ThreadID)
16 date.ranges = date.ranges[order(date.ranges$ThreadID),]
17
18 summary(date.ranges$Duration)
19 getmode(date.ranges$Duration)
20 boxplot(date.ranges$Duration)
21 hist(date.ranges$Duration)
22 t.test(date.ranges$Duration)
23
24 # Can we conclude that more popular threads (by post) are active for longer periods of
25 # time than less popular threads?
26 # No, we cannot conclude this at all. It's not significant.
27 top.posts = date.ranges[date.ranges$ThreadID %in% top.threads$ThreadID,]
28 bottom.posts = date.ranges[date.ranges$ThreadID %in% bottom.threads$ThreadID,]
29 t.test(top.posts$Duration, bottom.posts$Duration, alternative="greater")

```

Describe the summary statistics for the duration between new posts per thread. Perform t-test for the most popular and least popular threads.

```

1 thread.dur = as.data.frame(as.table(by(webforum$WC, list(webforum$Date, webforum$ThreadID), FUN=length)))
2 thread.dur = na.omit(thread.dur)
3 thread.dur = thread.dur[,-3]
4 colnames(thread.dur) = c("Date", "ThreadID")
5
6 # Convert to date data type.
7 thread.dur$Date = as.Date(thread.dur$Date, "%Y-%m-%d")
8 thread.dur$ThreadID = as.character(thread.dur$ThreadID)
9
10 avg.thread.dur = thread.dur %>% group_by(ThreadID) %>%
11   arrange(Date) %>%
12   summarize(Average = as.numeric(mean(diff(Date))))
13
14 avg.thread.dur = avg.thread.dur[order(avg.thread.dur$Average),]
15
16 summary(avg.thread.dur$Average)
17 getmode(avg.thread.dur$Average)
18 boxplot(avg.thread.dur$Average)
19 hist(avg.thread.dur$Average)
20 t.test(avg.thread.dur$Average)
21
22 # Can we conclude that more popular threads (by post) have shorter average durations
23 # between posts than less popular threads? No, not significant.
24 top.posts = avg.thread.dur[avg.thread.dur$ThreadID %in% top.threads$ThreadID,]
25 bottom.posts = avg.thread.dur[avg.thread.dur$ThreadID %in% bottom.threads$ThreadID,]
26 t.test(top.posts$Average, bottom.posts$Average, alternative="less")

```

Identify the main linguistic variables in the threads that have a lot of posts vs. threads that do not.

```

1 high.liwc = webforum[webforum$ThreadID %in% top.threads$ThreadID,]
2 low.liwc = webforum[webforum$ThreadID %in% bottom.threads$ThreadID,]
3
4 high.liwc = high.liwc[, 6:29]
5 low.liwc = low.liwc[, 6:29]
6
7 high.liwc = melt(high.liwc)

```

```

8 high.liwc.mean <- as.data.frame(as.table(by(high.liwc$value, list(high.liwc$variable),
9   mean)))
10 colnames(high.liwc.mean) = c("LIWC", "Proportion")
11 high.liwc.mean = cbind(high.liwc.mean, Class="Threads with Most Number of Posts")
12
13 low.liwc = melt(low.liwc)
14 low.liwc.mean <- as.data.frame(as.table(by(low.liwc$value, list(low.liwc$variable), mean
15   )))
16 colnames(low.liwc.mean) = c("LIWC", "Proportion")
17 low.liwc.mean = cbind(low.liwc.mean, Class="Threads with Lowest Number of Posts")
18
19 # Row bind
20 liwc.mean = rbind(high.liwc.mean, low.liwc.mean)
21
22 # Bar chart
23 g = ggplot(aes(fill=Class, y=Proportion, x=LIWC), data=liwc.mean)
24 g = g + geom_bar(position="dodge", stat="identity") + coord_flip()
25 g = g + theme_minimal()
26 g
27 ggsave("~/Desktop/FIT3152_resources/bar_chart.png", g, width = 20, height = 13, units =
28   "cm")
29 liwc.vars = as.vector(unique(high.liwc$variable))
30
31 # higher word count liwc variable proportion greater than lower word count
32 for (x in liwc.vars) {
33   print(x)
34   print(t.test(high.liwc$value[high.liwc$variable==x], low.liwc$value[low.liwc$variable
35     ==x], alternative = "greater"))
36 }
37
38 # higher word count liwc variable proportion greater than lower word count
39 for (x in liwc.vars) {
40   print(x)
41   print(t.test(high.liwc$value[high.liwc$variable==x], low.liwc$value[low.liwc$variable
42     ==x], alternative = "less"))
43 }
```

Analyse Thread 252620. The same analysis is done for threads 532649, 283958, 145223, 472752, 127115.

```

1 thread.details = webforum[webforum$ThreadID=="252620",]
2 # Group by MonthYear.
3 thread.details$MonthYear = as.Date(format(thread.details$date, "%Y-%m-01"))
4
5 liwc.thread = aggregate(thread.details[6:29], list(thread.details$date), FUN=mean)
6 colnames(liwc.thread)[1] = "Date"
7
8 melted.thread = melt(liwc.thread, id.vars=c("Date"))
9 colnames(melted.thread)[2] = "LIWC"
10 colnames(melted.thread)[3] = "Proportion"
11
12 liwc.thread.month = aggregate(thread.details[6:29], list(thread.details$MonthYear), FUN=
13   mean)
14 colnames(liwc.thread.month)[1] = "DateMonth"
15
16 melted.thread.month = melt(liwc.thread.month, id.vars=c("DateMonth"))
17 colnames(melted.thread.month)[2] = "LIWC"
18 colnames(melted.thread.month)[3] = "Proportion"
19
20 attach(melted.thread)
21 # plot faceted line graph
22 g = ggplot(melted.thread, aes(x=Date, y=Proportion)) + geom_jitter(colour="cadetblue3",
23   alpha=0.4)
```

```

22 g = g + facet_wrap(~LIWC, ncol=4)
23 g = g + geom_line(data=melted.thread.month, aes(x= DateMonth, y= Proportion), colour="black")
24 g
25
26 liwc.thread.analytic = aggregate(thread.details[6:9], list(thread.details$Date), FUN=
  mean)
27 colnames(liwc.thread.analytic)[1] = "Date"
28 melted.thread.analytic = melt(liwc.thread.analytic, id.vars=c("Date"))
29 colnames(melted.thread.analytic)[2] = "LIWC"
30 colnames(melted.thread.analytic)[3] = "Proportion"
31
32 # plot faceted line graph, analytic, clout, authentic and tone
33 g = ggplot(melted.thread.analytic, aes(x=Date, y=Proportion)) + geom_jitter(colour="cadetblue3", alpha=0.4)
34 g = g + facet_wrap(~LIWC, nrow=4)
35 g = g + stat_smooth(method = 'lm', aes(colour = 'Trend'), se = FALSE)
36 g = g + labs(color="Legend")
37 g
38
39 ggsave("~/Desktop/FIT3152_resources/LIWC_over_time1pt2.png", g, width = 20, height = 13,
  units = "cm")
40 detach(melted.thread)

```

5.3.3 Part 3

Analyse for March- the same analysis is done for April and May.

```

1 webforum$Year = as.Date(format(webforum$date, "%Y-%m-01"))
2 thread = webforum[webforum$Year=="2005-03-01",]
3
4 # Group by number of threads and authors.
5 author.threads = as.data.frame(as.table(by(thread$WC, list(thread$ThreadID, thread$AuthorID), FUN=length)))
6 colnames(author.threads) = c("ThreadID", "AuthorID", "NumPosts")
7 author.threads = na.omit(author.threads)
8 author.threads = author.threads[,-3]
9 author.threads = as.data.frame(as.table(by(author.threads$AuthorID, list(author.threads$ThreadID), FUN=length)))
10 colnames(author.threads) = c("ThreadID", "NumAuthors")
11 author.threads = author.threads[order(author.threads$NumAuthors),]
12
13 # There are 46 authors.
14 length(unique(thread$AuthorID))
15
16 # There are 7 threads in this month.
17 # Examine each of these threads.
18 thread.id = as.vector(author.threads$ThreadID)
19
20 threads = vector("list", length(thread.id))
21
22 for (i in seq_along(thread.id)) {
23   threads[[i]] = webforum[webforum$ThreadID==thread.id[i] & webforum$Year==as.Date(
24     format("2005-03-01")),]
25 }
26
27 author.id = vector()
28 # Examine the social network within each thread.
29 for (i in seq_along(thread.id)) {
30   author.id = c(author.id, unique(threads[[i]][2]))

```

```

30 }
31
32 # graph adjacency matrices
33 matrices = vector("list", length(author.id))
34 for (i in seq_along(author.id)) {
35   x = length(author.id[[i]])
36   matrices[[i]] = matrix(runif(x*x, min=1, max=1), ncol=x)
37   diag(matrices[[i]]) = 0
38   colnames(matrices[[i]]) = author.id[[i]]
39   rownames(matrices[[i]]) = author.id[[i]]
40 }
41
42 # graph social networks
43 graphs = vector("list", length(matrices))
44 par(mfrow=c(2, 3))
45 for (i in seq_along(matrices)) {
46   graphs[[i]] = graph.adjacency(matrices[[i]], mode="undirected", weighted=NULL)
47   if (i > 2) {
48     par(mar=c(0,0,0,0)); plot(graphs[[i]], layout=layout.circle, vertex.color="red",
49       vertex.label.cex = c(0.5),
50       vertex.label.color = "black",
51       vertex.frame.color = "white",
52       vertex.size=20)
53   }
54 }
55
56 par(mfrow=c(1, 1))
57 # merge graphs into huge social network
58 mar = graph_from_literal()
59 for (i in seq_along(graphs)) {
60   mar = (mar %u% graphs[[i]])
61 }
62
63 V(mar)$size <- log(strength(mar)) * 3 + 3
64 V(mar)$label <- ifelse( strength(mar)>=10, V(mar)$name, NA )
65 par(mar=c(0,0,0,0))
66 png('march_graph.png', width = 2000, height = 1700)
67 plot(mar, layout=layout.fruchterman.reingold,
68       vertex.frame.color="white",
69       vertex.label.color="black",
70       vertex.label.cex=c(1.5))
71 dev.off()
72
73 # Calculate statistics
74 mar
75
76 # Degree Distribution
77 hist(degree(mar), breaks = 30, col="grey")
78 deg = as.table(degree(mar))
79
80 # Betweenness
81 bet = (format(as.table(betweenness(mar)), digits=2))
82
83 # Closeness
84 clo = format(as.table(closeness(mar)), digits=2)
85
86 # Eigenvector
87 eig = as.table(evcent(mar)$vector)
88
89 # Put into data frame
90 summary = as.data.frame(cbind(deg, bet, clo, eig))
91 summary

```

```

92 # Diameter
93 diameter(mar)
94
95 # Average Path Length
96 average.path.length(mar)
97
98 # Cliques
99 # table(sapply(cliques(mar), length))
100
101 attach(summary)
102 deg = as.numeric(deg)
103 bet = as.numeric(bet)
104 clo = as.numeric(clo)
105 eig = as.numeric(eig)
106
107 hdeg = summary[order(-deg),]
108 head(hdeg)
109
110 hclo = summary[order(-clo),]
111 head(hclo)
112
113 hbet = summary[order(-bet),]
114 head(hbet)
115
116 heig = summary[order(-eig),]
117 head(heig)
118
119 detach(summary)
120
```

Analyse for March and April. The same analysis is done to combine March, April and May.

```

1 comb = (mar %u% apr)
2
3 V(comb)$size <- log(strength(comb)) * 3 + 3
4 V(comb)$label <- ifelse( strength(comb)>=10, V(comb)$name, NA )
5 par(mar=c(0,0,0,0))
6 png('march_april_graph.png', width = 2000, height = 1500)
7 plot(comb, layout=layout.fruchterman.reingold,
8       vertex.frame.color="white",
9       vertex.label.color="black",
10      vertex.label.cex=c(1.3))
11 dev.off()
12
13 # Calculate statistics
14 comb
15
16 # Degree Distribution
17 hist(degree(comb), breaks = 30, col="grey")
18 deg = as.table(degree(comb))
19
20 # Betweenness
21 bet = (format(as.table(betweenness(comb)), digits=2))
22
23 # Closeness
24 clo = format(as.table(closeness(comb)), digits=2)
25
26 # Eigenvector
27 eig = as.table(evcent(comb)$vector)
28
29 # Put into data frame
```

```
30 summary = as.data.frame(cbind(deg, bet, clo, eig))
31 summary
32
33 # Diameter
34 diameter(comb)
35
36 # Average Path Length
37 average.path.length(comb)
38
39 # Cliques
40 # table(sapply(cliques(comb), length))
41
42 attach(summary)
43 deg = as.numeric(deg)
44 bet = as.numeric(bet)
45 clo = as.numeric(clo)
46 eig = as.numeric(eig)
47
48 hdeg = summary[order(-deg),]
49 head(hdeg)
50
51 hclo = summary[order(-clo),]
52 head(hclo)
53
54 hbet = summary[order(-bet),]
55 head(hbet)
56
57 heig = summary[order(-eig),]
58 head(heig)
59
60 detach(summary)
```