

BUSINESS INTELLIGENCE AND DATA MINING (BIDB17-4)

TERM 4

***“Assessing Consumer Financial Complaints-
Driving insights using Sentiment Analysis”***

GROUP B16

<i>Ritwika Majumdar</i>	<i>B17162</i>
<i>Sreela Basu Roy</i>	<i>B17174</i>
<i>Priyadarshi Dasgupta</i>	<i>H17155</i>
<i>Sanchit Agarwal</i>	<i>H17163</i>

**Business Management
XLRI Jamshedpur**

Contents

1. Data Preparation	3
2. Sentiment Analysis.....	3
1.1. Observations	5
1.2. Business Insights	15
2. Frequency Analysis.....	15
2.1 Observations	15
2.2 Business Insights	19
3. Decision Tree using RPart	19

INTRODUCTION

In this project we are analysing the customer complaints for various financial services such as Credit card or prepaid card, mortgage, vehicle loan or lease, debt collection, checking or savings account and so on received in 2018 in USA. The dataset was obtained from the website (<https://catalog.data.gov/>) which has public records made available.

Data Source: <https://catalog.data.gov/dataset/consumer-complaint-database>

The data acquired was cleaned to arrive in the format:

<i>Date received</i>	The date that the complaint was received on
<i>Product</i>	The type of product/service the customer identified
<i>Sub-product</i>	The type of sub product identified by the customer
<i>Issue</i>	The issue identified by the customer in the complaint
<i>Sub-issue</i>	The sub issue identified by the customer in the complaint
<i>Consumer complaint narrative</i>	Complaint recorded verbatim in text format
<i>Company public response</i>	The company's public response to the consumer complaint
<i>Company</i>	The company the complaint is against
<i>State</i>	Customer's state as per mailing address
<i>ZIP code</i>	Customer's ZIP code as per mailing address
<i>Tags</i>	Data that supports easier searching on behalf of consumers
<i>Consumer consent provided?</i>	If the consumer consents to publishing the complaint narrative
<i>Submitted via</i>	Medium of submission of the complaint
<i>Date sent to company</i>	The date that the complaint is forwarded to the company
<i>Company response to consumer</i>	How the company responded to the complaint
<i>Timely response?</i>	If the response by the company was given in appropriate time
<i>Complaint ID</i>	Unique identification number of a complaint

We will be doing sentiment analysis on the consumer complaints narrative to find out about the specific nature of complaints via text mining. Since these are complaints, most/all of the sentiments will be negative. However, we shall look at the magnitudes across categories to understand the dissatisfaction. We will also be looking at the frequencies of complaints based on different states and across different product categories to understand the major problem areas that a financial service provider faces.

METHODOLOGY

1. Data Preparation

The dataset received from the site required data cleaning and preparation before any analysis could be performed. The data was filtered for 2018 data, and irrelevant columns were removed. Rows with corrupted or blanks cells were also removed, thus bringing down 5 Lakh records to around 1 Lakh.

2. Sentiment Analysis

Sentiment analysis is the process of identifying and categorizing opinions expressed in a piece of text to determine broadly what the author's attitude towards that particular subject is. Lots of major companies employ algorithms to translate consumer feedback into sentiments for gaining insight into a consumer's thought process. This approach helps them to understand what the consumer expects and how product/service matches up to those expectations.

Generally, it is done as a binary distinction between positive and negative sentiments, but it may also be more fine-tuned to more specific sentiments felt by the author such as anger, fear, joy, etc. Of course, it is not always a 100% accurate, but it is a useful way to measure the emotions portrayed in large blocks of text.

The following libraries were loaded to help us perform this analysis:

```
library(tidyverse)
library(tidytext)
library(glue)
library(stringr)
```

The tidytext package performs sentiment analysis using the following method:

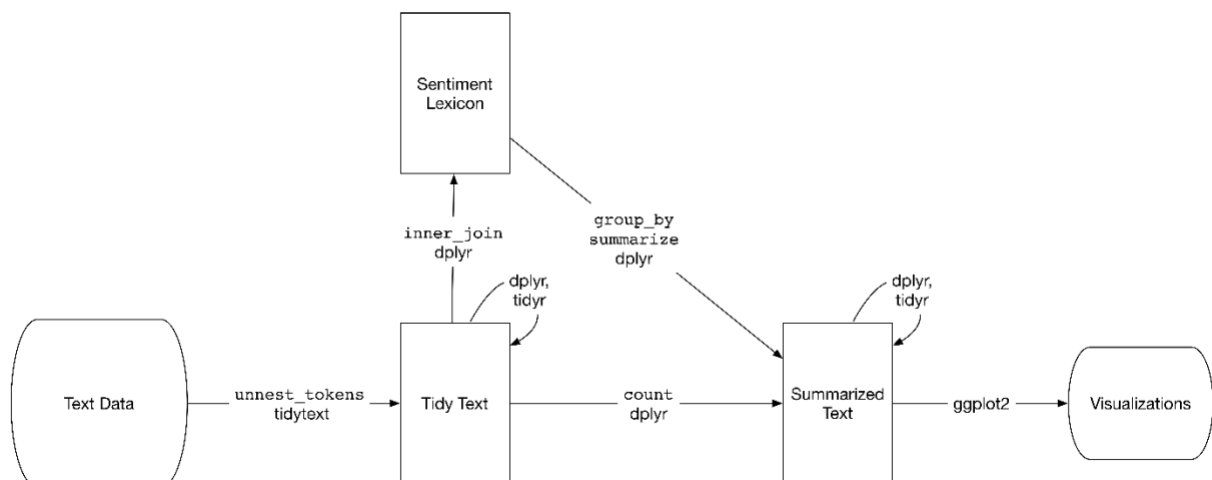


Figure: A flowchart of a typical text analysis that uses tidytext for sentiment analysis.

The code used on R to perform analysis for the same is as follows:

```
#Sentiment Analysis
#Extract complaint narrative
dataset_main1<-dataset_main[!(dataset_main$consumer.complaint.narrative==""),]

#Find unique product types
product = unique(dataset_main1[c("Product")])

#Clean data for sentiment analysis and write complaints for each product to separate text file
cleanData = function(complaint, fileName){
  complaint = gsub("[[:punct:]]", "", complaint)
  complaint = gsub("@\\w+", "", complaint)
  complaint = gsub("[[:digit:]]", "", complaint)
  complaint = gsub(" ", "", complaint)
  complaint = gsub("[ \\t]{2,}", " ", complaint)
  complaint = gsub("xx", " ", complaint)
  complaint = gsub("xxx", " ", complaint)
  complaint = gsub("xxxx", "", complaint)
  complaint = gsub("xxxxxxxx", "", complaint)
  complaint <- tolower(complaint)
  complaint <- complaint[!is.na(complaint)]
  complaint <- gsub("\\s+", " ", complaint)
  complaint <- trimws(complaint)
  filePath = paste("input/", fileName, ".txt", sep="")
  write.table(complaint, file=filePath, sep = "\t", row.names = FALSE)
}

#Loop through products and call cleanData function
for(i in 1:9){
  complaint= dataset_main2[,6]
  cleanData(complaint,product[i,1])
  # cleanData(sentiment_dataset,product[i,1])
}

#Store list of filenames (complaint files)
files <- list.files("c://Users/sree1/OneDrive/XLRI/Term 4/BIDM/Project/input")
```

```
# write a function that takes the name of a file and returns the # of positive
# sentiment words, negative sentiment words, and the difference
GetSentiment <- function(file){
  # get the file
  fileName <- glue("input/", file, sep = "")
  # get rid of any sneaky trailing spaces
  fileName <- trimws(fileName)

  df1 <- read.table(fileName, header=TRUE, fill = TRUE)

  # read in the new file
  fileText <- glue(read_file(fileName))
  # remove any dollar signs (they're special characters in R)
  fileText <- gsub("\\$", "", fileText)

  # tokenize
  tokens <- data_frame(text = fileText) %>% unnest_tokens(word, text)

  # get the sentiment from the first text:
  sentiment <- tokens %>%
    inner_join(get_sentiments("bing")) %>% # pull out only sentiment words
    count(sentiment) %>% # count the # of positive & negative words
    spread(sentiment, n, fill = 0) %>% # made data wide rather than narrow
    mutate(sentiment = positive - negative) %>% # # of positive words - # of negative words
    mutate(file = file) %>% # add the name of our file
    mutate(product = file) # add product

  # return our sentiment dataframe
  return(sentiment)
}

# file to put our output in
sentiments <- data_frame()

# get the sentiments for each file in our dataset
for(i in files){
  sentiments <- rbind(sentiments, GetSentiment(i))
}

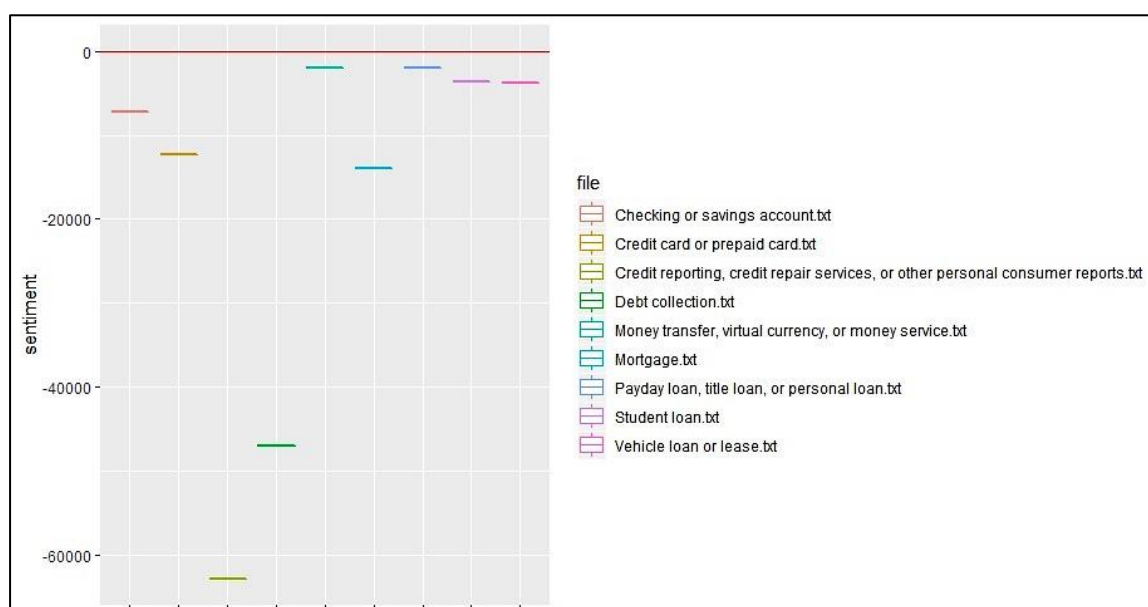
#plot overall sentiment for each product type
ggplot(sentiments, aes(x = file, y = sentiment, color = file)) +
  geom_hline(yintercept = 0, color = "red") +
  theme(axis.text.x = element_blank()) +
  geom_boxplot() # draw a boxplot for each product
```

1.1. Observations

The **consumer complaint narratives** for each product type were analysed as follows:

- 1.1.1. Sentiment Analysis-** This is the classification by emotion, each complaint is classified by emotion and this is the histogram representing each group of emotion, using the bing lexicon (part of the tidytext package)
- 1.1.2. Word Cloud-** A word cloud represents all the words available in the database, the size of each word tells its frequency in the database, the bigger the word is, the more frequent that word appears in the database.

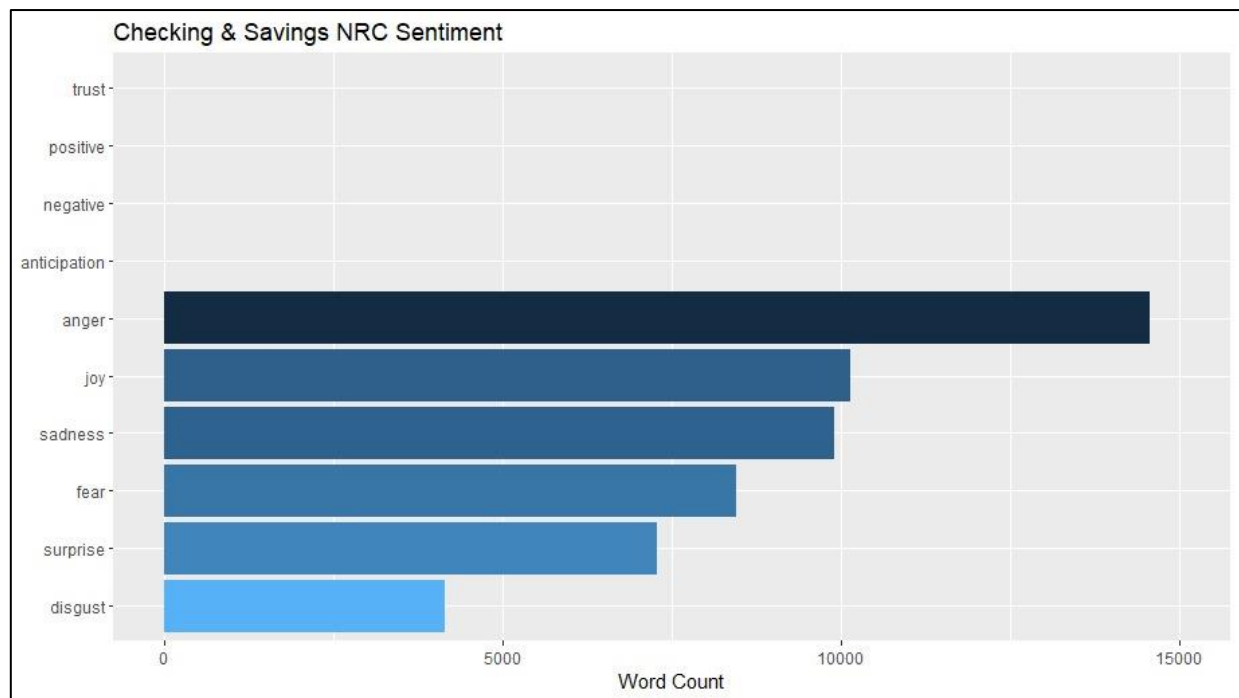
Furthermore, we did sentiment analysis and assessed word clouds on the complaints for each different type of financial service that is provided.



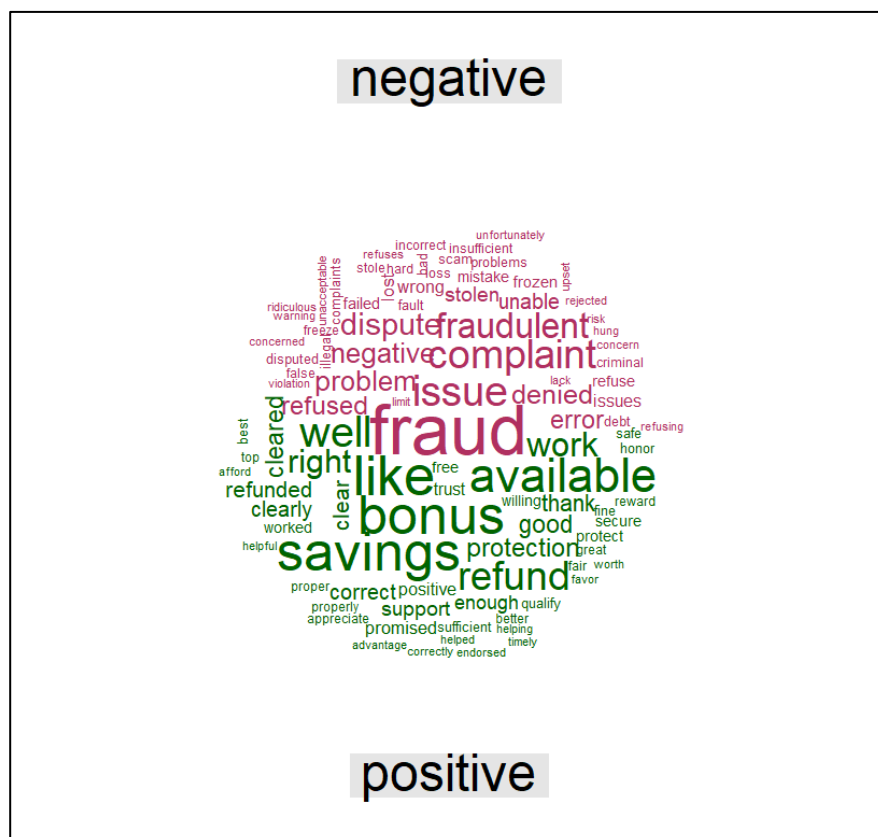
The sentiment analysis of the consumer complaint narratives using the **tidytext** packages shows that all the narratives are primarily negative, with Credit Reporting, credit repair services or other personal consumer reports being most negative. **The NRC lexicon** (part of the **tidytext** package) was used to classify the words in the narratives into positive and negative, and an overall score for each product type was arrived at.

Following are the sentiment analysis histograms and word clouds for complaints on each financial service:

A. For checking and savings accounts:

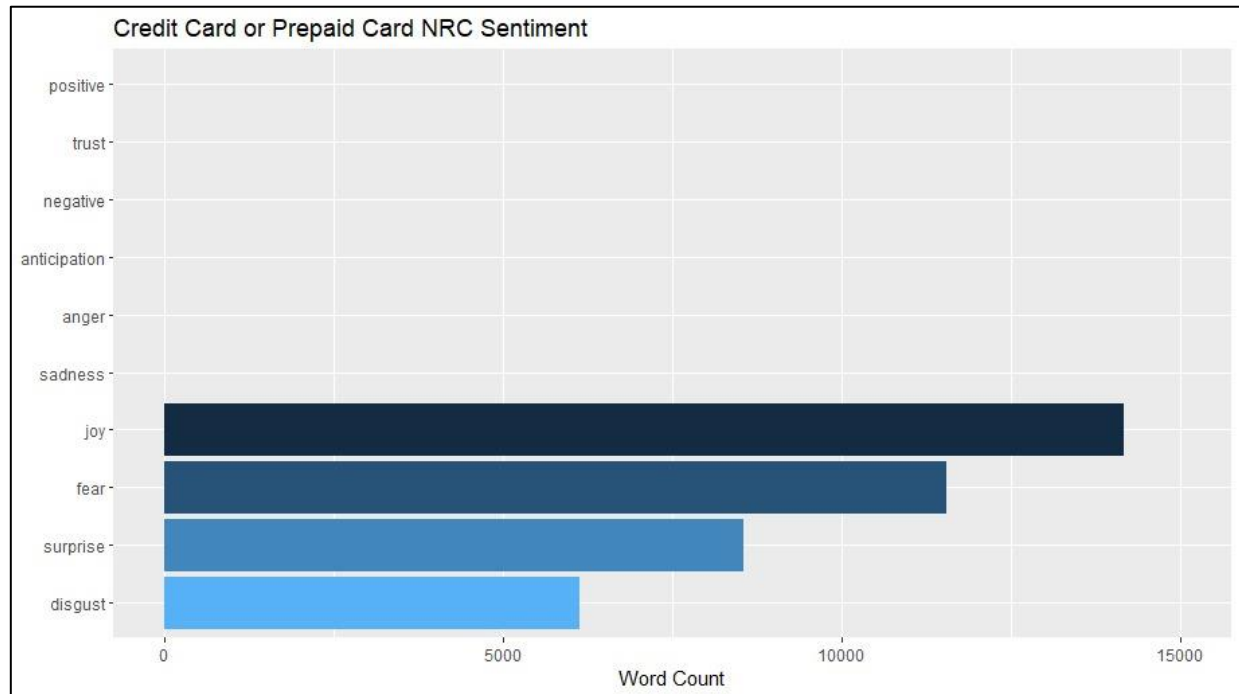


Here, anger is the emotion with most complaints followed by joy, sadness, fear and so on.

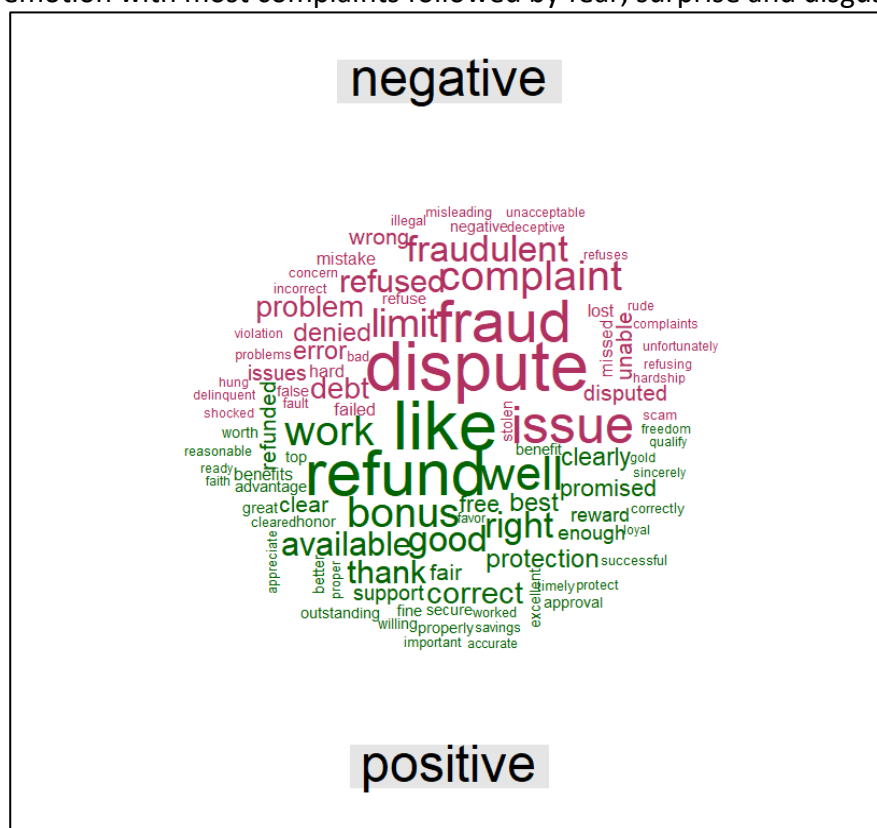


“Fraud” is the word with the highest frequency in our given database, with respect to the Checking and Savings Accounts.

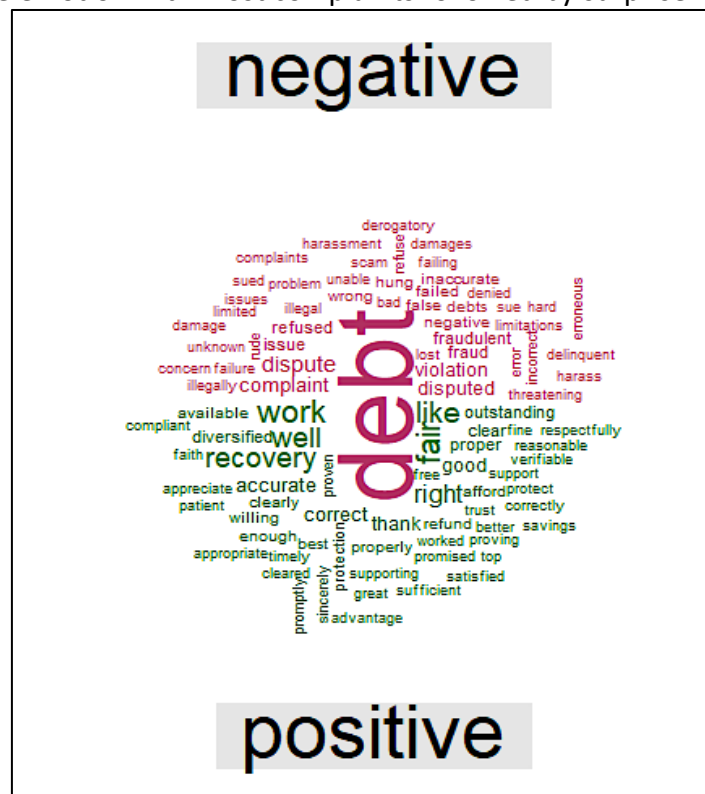
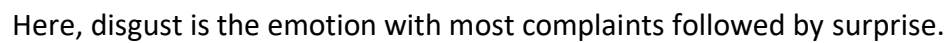
B. For credit card or prepaid cards:



Here, joy is the emotion with most complaints followed by fear, surprise and disgust.

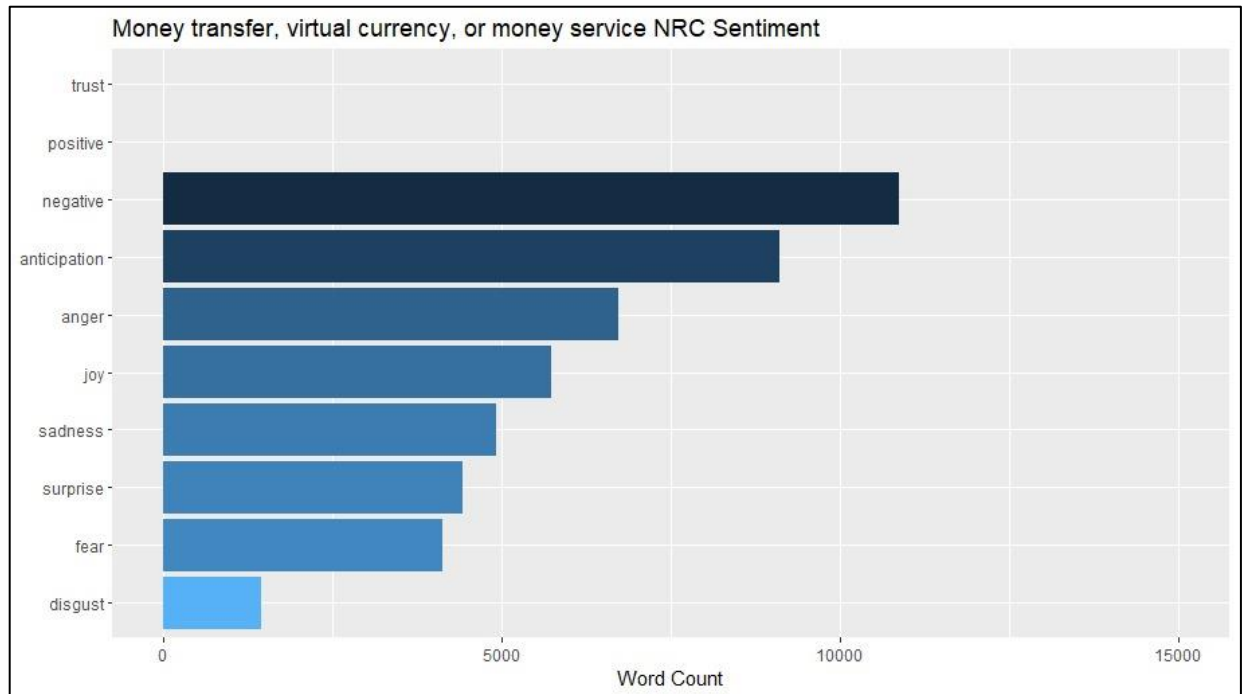


C. For debt collection:

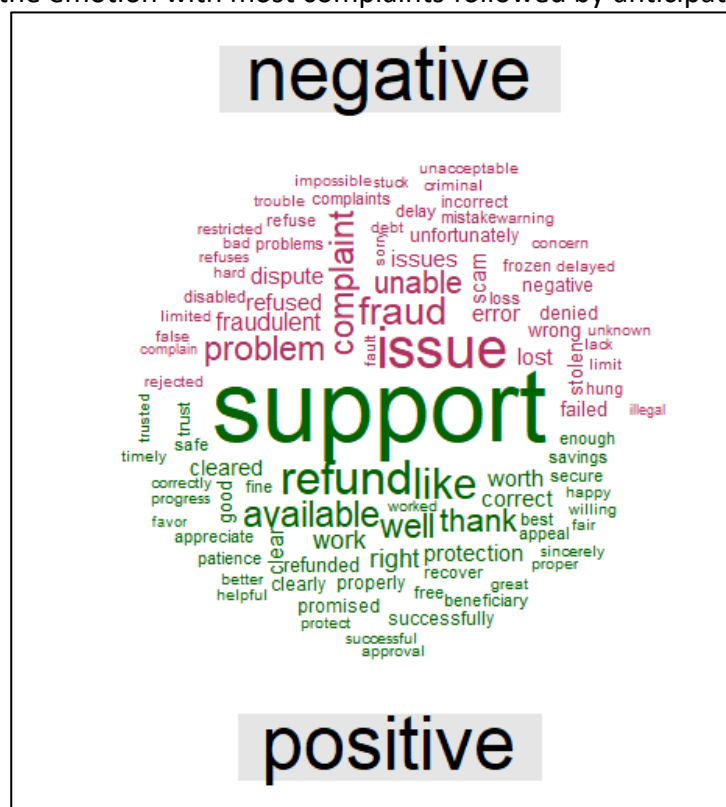


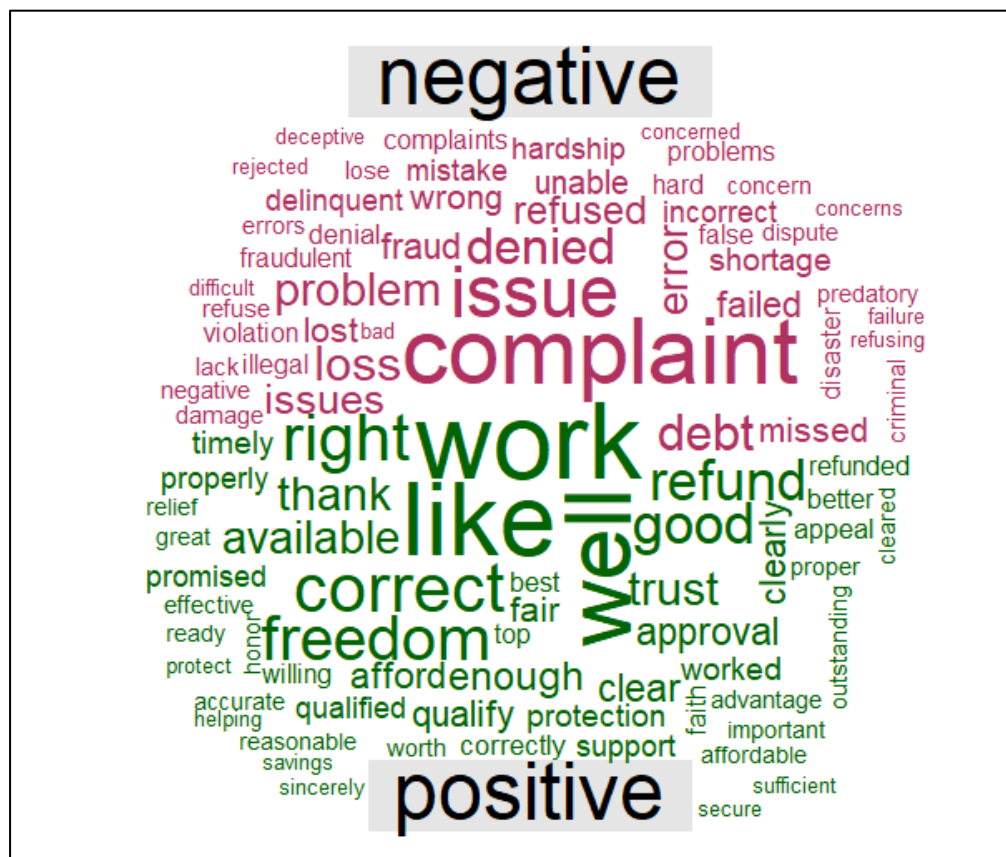
“Debt” is the word with the highest frequency in our given database, with respect to the Debt Collection.

D. For money transfer, vital currency or money services:



Here, negative is the emotion with most complaints followed by anticipation, anger and so on.





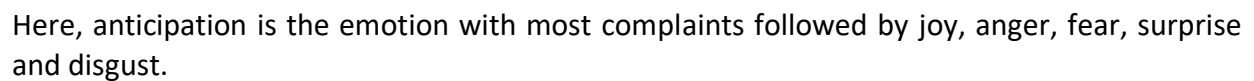
negative

refusing delinquent problems misleading deceptive complaints unfortunately impossible harassment hardship lost negative dispute problem penalty concerns refused missed fraudulent threat debt issue complaint unable error violations wrong denied poor predatory illegal issues hard struggling lied hung disputed scam revoked excessive fraud refuse properly clearly timely available correctly clear promised qualify prosper advantage fine right enough speedy support reasonable appreciate savings accurate helping worked easy afford good top outstanding recovery willing refunded great protection grace better honor thank trust approval proper fast

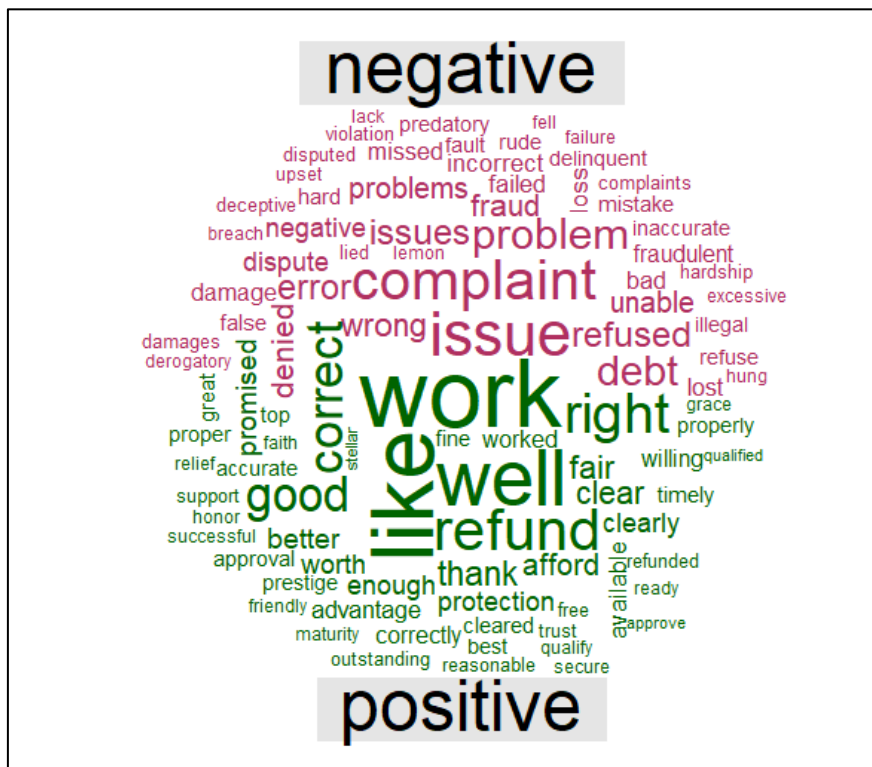
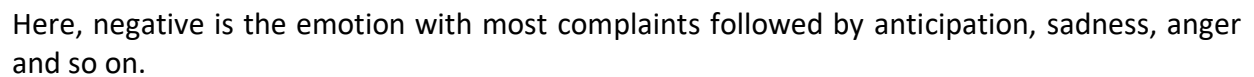
work like well right best correct refund fair support

positive

G. For student loan:



H. For vehicle loan or lease:



“Work”, “Like”, “Refund” are some of the words with the highest frequency in our given database, with respect to Vehicle Loan or Lease.

```
#plot nrc plot
undesirable_words <- c("dear","sir","madam")

#Store list of filenames (complaint files)
files <- list.files("C://Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project/input")

#glue path to file name
fileName <- glue("C://Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project/input/", files[9], sep = "")

#Remove whitespace
fileName <- trimws(fileName)

fileText <- glue(read_file(fileName))

tokens <- data_frame(text = fileText) %>% unnest_tokens(word, text)
tokens %>%
  filter(!nchar(word) < 3) %>% #remove words with less than 3 characters
  filter(!word %in% undesirable_words) %>%
  anti_join(stop_words)

savings_nrc <- tokens %>%
  inner_join(get_sentiments("nrc"))

nrc_plot <- savings_nrc %>%
  group_by(sentiment) %>%
  summarise(word_count = n()) %>%
  ungroup() %>%
  mutate(sentiment = reorder(sentiment, word_count)) %>%
  #Use 'fill = -word_count' to make the larger bars darker
  ggplot(aes(sentiment, word_count, fill = -word_count)) +
  geom_col() +
  guides(fill = FALSE) + #Turn off the legend
  labs(x = NULL, y = "word count") +
  scale_y_continuous(limits = c(0, 15000)) + #Hard code the axis limit
  ggtitle("vehicle loan or lease NRC Sentiment") +
  coord_flip()

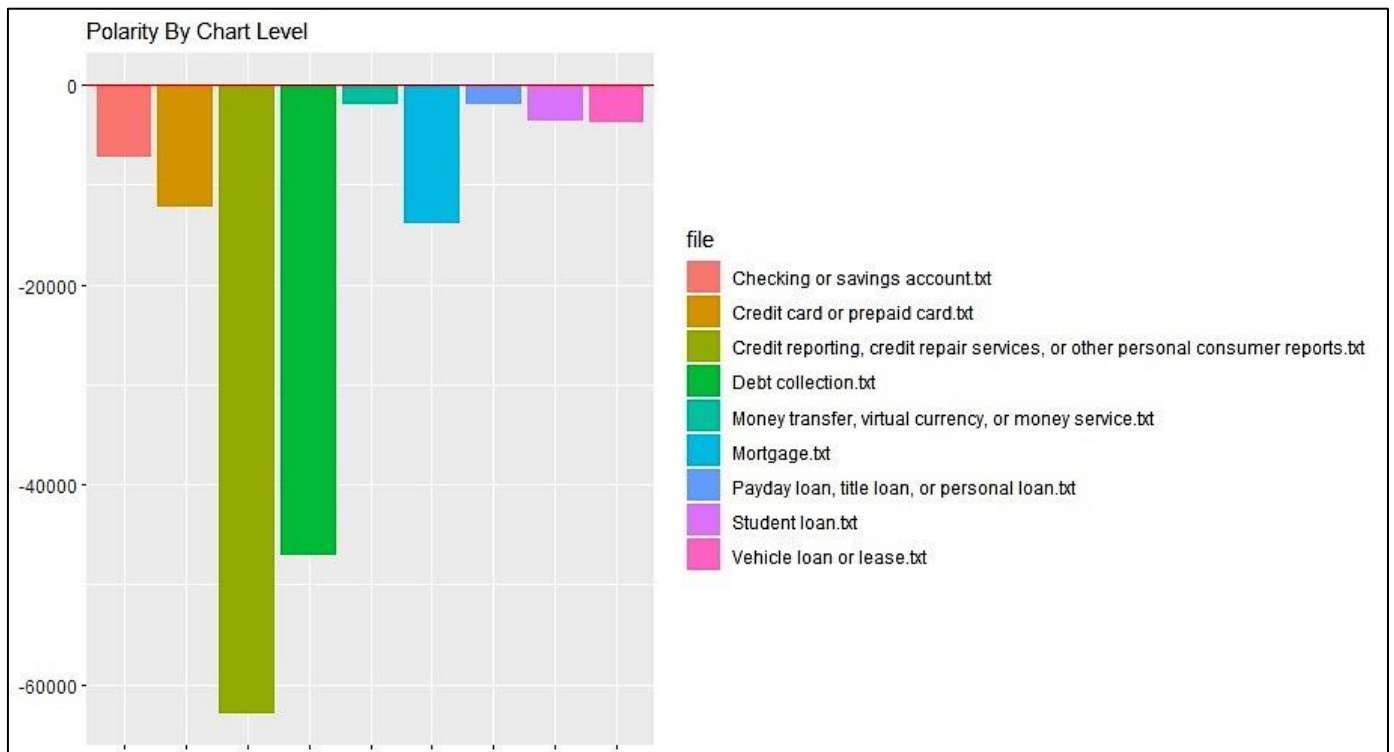
plot(nrc_plot)
```

1.1.3. Polarity-

Additionally, we found polarity of the sentiments used in all these categories with results as in the table below:

```
> sentiments
# A tibble: 9 x 5
  negative positive sentiment file
  <dbl>      <dbl>      <dbl> <chr>
1    20214    12994    -7220 checking or savings account.txt
2    32347    20063   -12284 credit card or prepaid card.txt
3   111110    48239   -62871 credit reporting, credit repair services, or ot~
4    73950    26913   -47037 Debt collection.txt
5     8994     7006   -1988 Money transfer, virtual currency, or money serv~
6    40257    26318  -13939 Mortgage.txt
7     5221     3270   -1951 Payday loan, title loan, or personal loan.txt
8    12933     9360   -3573 student loan.txt
9     9011     5294   -3717 vehicle loan or lease.txt
```

This led to the following polarity comparison between the categories of complaints:



Through the polarity analysis using bing lexicon in tidytext, we can infer that Credit Reporting, credit repair services or other personal consumer reports cause the maximum negative complaints, followed by Debt collection & Mortgage

1.2. Business Insights

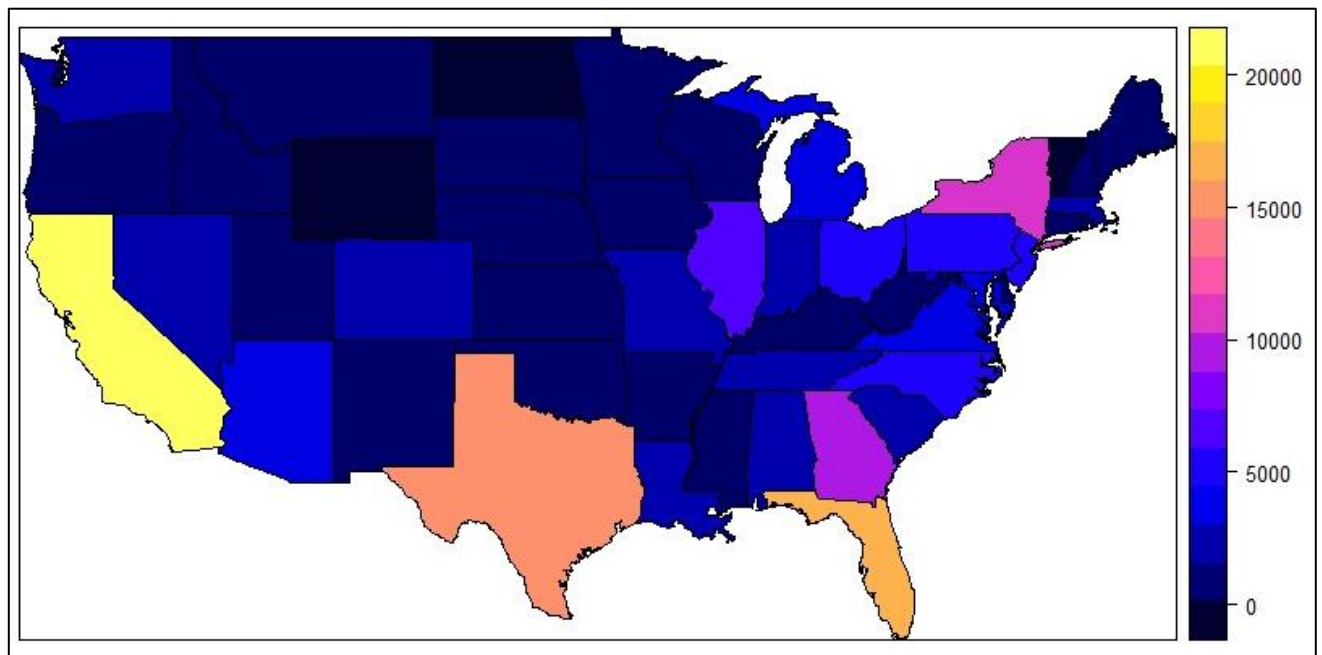
1. **Credit Reporting, credit repair services, or other personal consumer reports and Debt Collection** create the largest number of complaints amongst the users. Companies should focus on these major criteria so that they are able to track and handle issues in an effective manner.
2. From the word cloud and sentiment analysis, "**Fraud**" is the leading opinion that runs amongst the public. So it is imperative to address these concerns in order to better reach the people and solve their problems better.

2. Frequency Analysis

The frequency of the different types of complaints according to various different kinds of parameters might help us understand problem areas. The frequency of different type of products would also help us understand which products are more problematic in general. We will be able to see what kinds of complaints arise in which areas.

2.1 Observations

Total heat map of the complaints received across the 50 states:

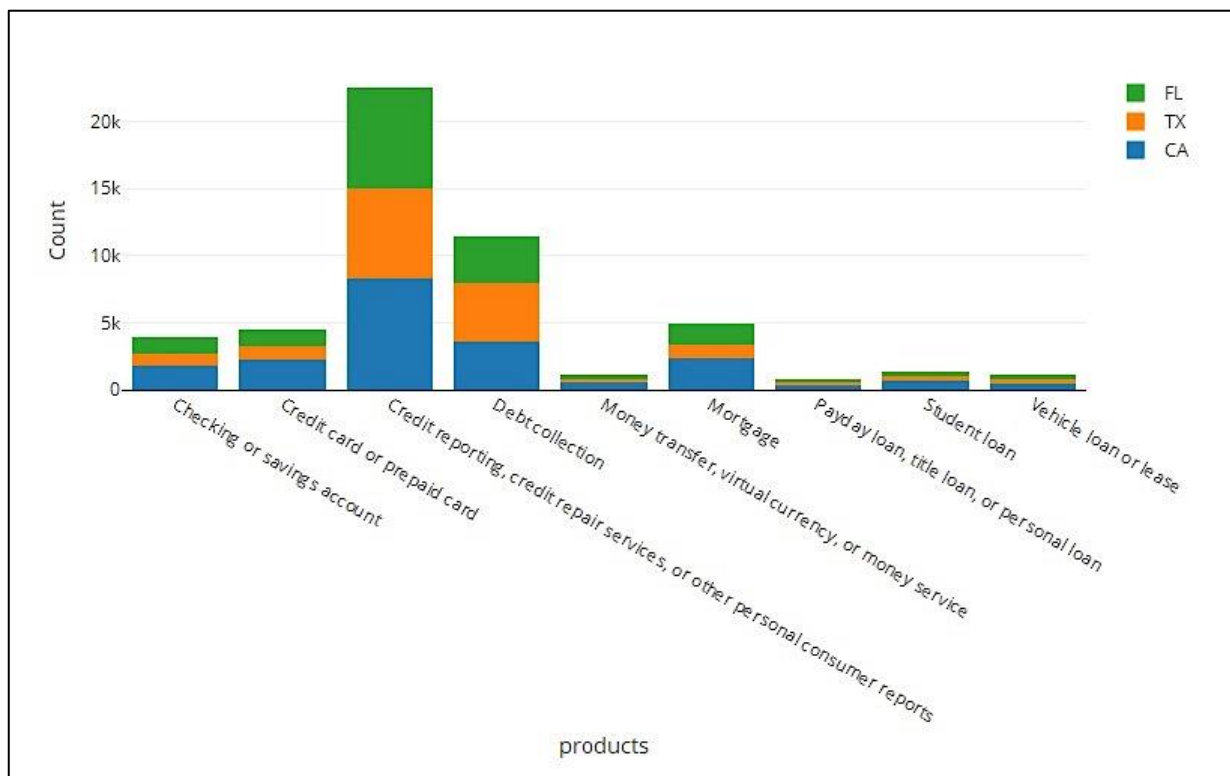


From this, we notice that the top 3 states in terms of complaints are:

- **Florida (FL)**
- **Texas (TX)**
- **California (CA)**

Taking a look at the break-up of the complaints in this area, we see:

For all types of financial complaints, California has the highest number of complaints. Overall, Credit Reporting, credit repair services, or other personal consumer reports have registered the highest complaints across all the three states. Whereas, Payday loan, title loan, or personal loan has the least amount of complaints across these regions of states.



```
#plot complaint frequencies on the US map
library(maps)
library(maptools)
library(sp)
library(plotly)
b<-as.data.frame(table(dataset_main$State))
str(b)
names(b)[1] <- 'state.abb'
str(b)
b$states <- tolower(state.name[match(b$state.abb, state.abb)])
str(b$states)
mapUSA <- map('state', fill = TRUE, plot = FALSE)
nms <- sapply(strsplit(mapUSA$names, ':'), function(x)x[1])
USApolygons <- map2SpatialPolygons(mapUSA, IDs = nms, CRS('+proj=longlat'))

idx <- match(unique(nms), b$states)
dat2 <- data.frame(value = b$Freq[idx], state = unique(nms))
row.names(dat2) <- unique(nms)

USAsp <- SpatialPolygonsDataFrame(USApolygons, data = dat2)
spplot(USAsp['value'])
```

```

#plot complaint frequencies by location
ca_data=subset(dataset_main, State == "CA")
ca_prod<-as.data.frame(table(ca_data$Product))
ca_prod_freq = ca_prod$Freq

tx_data=subset(dataset_main, State == "TX")
tx_prod<-as.data.frame(table(tx_data$Product))
tx_prod_freq = tx_prod$Freq

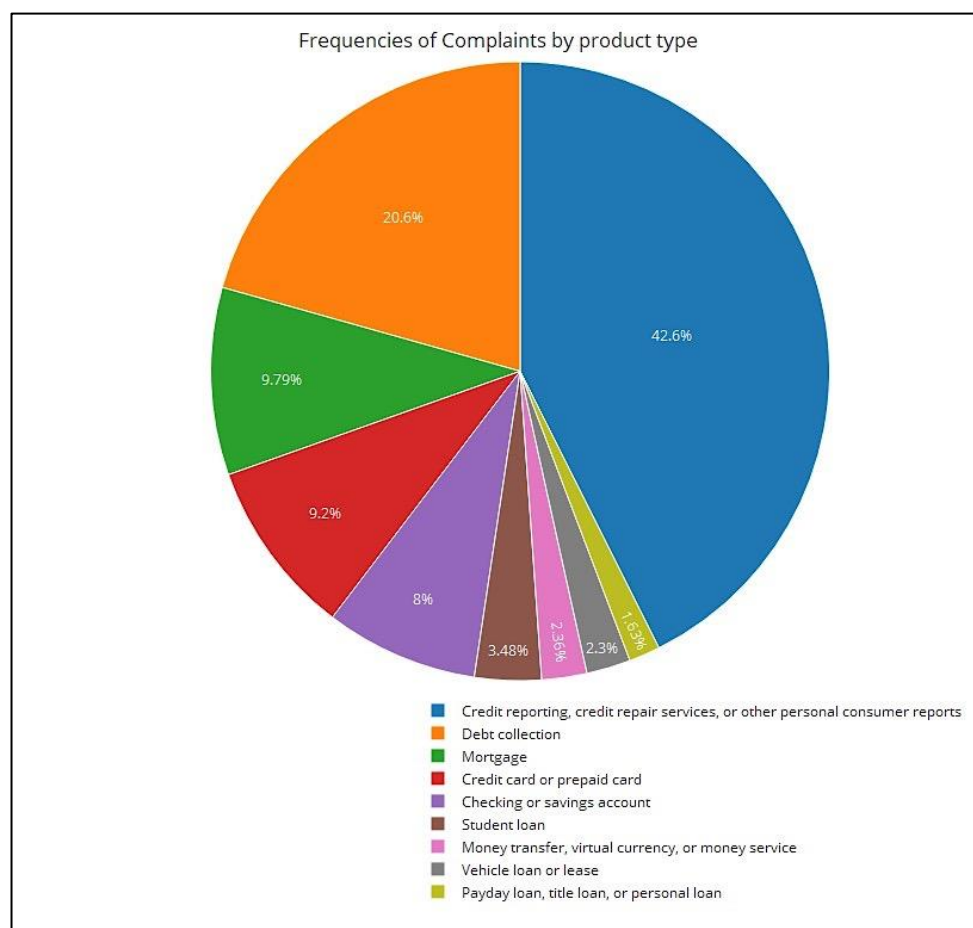
fl_data=subset(dataset_main, State == "FL")
fl_prod<-as.data.frame(table(fl_data$Product))
fl_prod_freq = fl_prod$Freq

products = ca_prod$Var1
data = data.frame(products,ca_prod_freq,tx_prod_freq)

p <- plot_ly(data, x = ~products, y = ~ca_prod_freq, type = 'bar', name = 'CA') %>%
  add_trace(y = ~tx_prod_freq, name = 'TX') %>%
  add_trace(y = ~fl_prod_freq, name = 'FL') %>%
  layout(yaxis = list(title = 'count'), barmode = 'stack')
p

```

Overall the break-up of complaints across the product categories are:

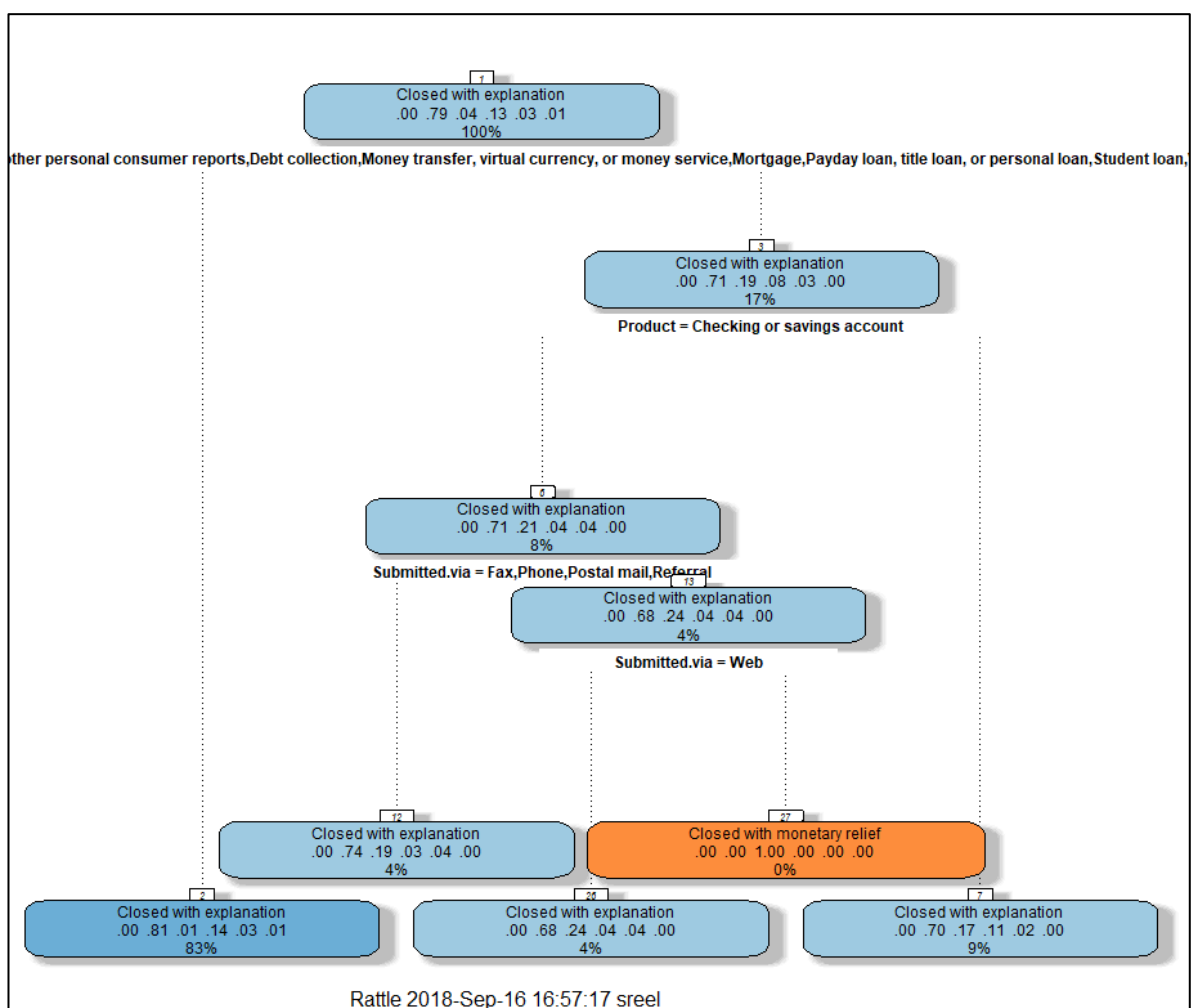


From the above pie chart we can infer that maximum number of complaints arise from the Credit reporting, credit repair services, or other personal consumer reports category, followed by debt collection and mortgage.

2.2 Business Insights

1. Considering the top 3 states with the highest number of complaints are Florida, Texas and California, companies need to be more proactive in these regions, especially if they are looking to open up new branches or areas of service.
2. Amongst these regions, special attention should be given to Credit Reporting, credit repair services, or other personal consumer reports and Debt Collection, which contribute to the highest portion of complaints.

3. Decision Tree using RPart



The decision tree classifies the consumer complaints on the basis of product type, method of submission to arrive at how the consumer responded to customers about their complaints.

The unique product types are as follows:

- ```
[1] "Checking or savings account "
[2] "Credit card or prepaid card "
```

[3] "Credit reporting, credit repair services, or other personal consumer reports "

[4] "Debt collection "

[5] "Money transfer, virtual currency, or money service "

[6] "Mortgage "

[7] "Payday loan, title loan, or personal loan "

[8] "Student loan "

[9] "Vehicle loan or lease "

- For **Products 3-9**, the 83% responses from the company are "Closed with Explanation".
- For **Product 2**, 9% responses from the company are "Closed with Explanation".
- For **Product 1**, 8% responses from the company are "Closed with Explanation". The breakup of this 8% is obtained from the following sub-cases:
  - 4% responses from the company are "Closed with Explanation"- those submitted via Fax, Phone, Postal Mail and Referral
  - 4% responses from the company are "Closed with Explanation"- those submitted via Web
- For **Product 1**, those which are submitted via e-mail have negligible (~0%) cases where monetary relief is given.

```
#decision tree
library(dplyr)
dataset_main3 <- dataset_main
decision_data = select(dataset_main3, Product, Submitted.via, Company.response.to.consumer)
str(decision_data)

library(rpart)
library(rattle)
ac<-rpart(Company.response.to.consumer~.,data=decision_data,method="class",control =rpart.control(minsplit =1,minbucket=1, cp=0))
summary(ac)
plot(ac, uniform=TRUE, main="Classification tree for Company Response to consumers")
text(ac,pretty=0)

fancyRpartPlot(ac)
```

# REFERENCES

1. <https://www.tidytextmining.com/sentiment.html>
2. <https://stackoverflow.com>

# CODES

```
1 library(tm)
2 library(RColorBrewer)
3 library(NLP)
4 library(wordcloud)
5 library(SnowballC)
6 library(plyr)
7 library(stringr)
8 library(Rstem)
9 library(tidyverse)
10 library(tidytext)
11 library(glue)
12 library(reshape2)
13 library(dplyr)
14 setwd('C:/Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project')
15
16 dataset_main = read.csv('2018.csv',header = T)
17
18 issue = dataset_main[,4]
19 complaint = dataset_main[,6]
20 # issue is taking data only for column 4 which has the issue in complaints.
21 # We have clustered only issue but any other column can be clustered as well.
22 ab=Corpus(VectorSource(issue))
23 #inspect(ab)
24
25 ac= tm_map(ab,tolower)
26 ac = tm_map(ac, removeWords, stopwords("english"))
27 ac = tm_map(ac, removePunctuation)
28 ac = tm_map(ac, removeNumbers)
29 ac = tm_map(ac, stripWhitespace)
30
31 ac = tm_map(ac, stemDocument, language = "english")
32 #statements to prepare the document
33
34 wordcloud(ac,min.freq = 4)
35 dtm = DocumentTermMatrix(ac)
36 #head(inspect(dtm1))
37 dtm1 = as.matrix(dtm)
38 #dtm1
39 kk = colSums(dtm1)
40 kk1 = sort(kk, decreasing = TRUE)
41 #head(kk1)
42
43 wtt = weightTfIdf(dtm)
44 #inspect(wtt)
45 wtt1 = as.matrix(wtt)
```

```

46 #wtt1
47
48
49 f1= function(x)
50 {
51 (sum(x^2))^0.5
52 }
53 f2 = function(y)
54 {
55 wtt1/apply(y, 1, f1)
56 }
57 am = f2(wtt1)
58
59 abc = 10
60
61 # for (i in 2:10)
62 # {
63 # abc[i]<-kmeans(am, centers = i, nstart = 10, iter.max = 10)$tot.withinss
64 # }
65 # plot(1:10,abc,type = "b")
66 ak = kmeans(am,5)
67 #ak
68
69 dd<-cbind(dataset_main,ak$cluster)
70 head(dd)
71
72
73
74
75
76
77 #Sentiment Analysis
78 #Extract complaint narrative
79 dataset_main1<-dataset_main[!(dataset_main$Consumer.complaint.narrative==""),]
80
81 #Find unique product types
82 product = unique(dataset_main1[c("Product")])
83
84 #Clean data for sentiment analysis and write complaints for each product to separate text file
85 CleanData = function(complaint, fileName){
86 complaint = gsub("[[:punct:]]", "", complaint)
87 complaint = gsub("@\\w+", "", complaint)
88 complaint = gsub("[[:digit:]]", "", complaint)
89 complaint = gsub(" ", " ", complaint)

```



```

90 complaint = gsub("[\t]{2,}", " ", complaint)
91 complaint = gsub("XX", " ", complaint)
92 complaint = gsub("XXX", " ", complaint)
93 complaint = gsub("XXXX", "", complaint)
94 complaint = gsub("XXXXXXXX", "", complaint)
95 complaint <- tolower(complaint)
96 complaint <- complaint[!is.na(complaint)]
97 complaint <- gsub("\\s+", " ", complaint)
98 complaint <- trimws(complaint)
99 filePath = paste("input/", fileName, ".txt", sep = "")
100 write.table(complaint, file=filePath, sep = "\t", row.names = FALSE)
101 }
102
103 #Loop through products and call CleanData function
104 for(i in 1:9){
105 complaint= dataset_main2[,6]
106 CleanData(complaint,product[i,1])
107 # CleanData(sentiment_dataset,product[i,1])
108 }
109
110 #Store list of filenames (complaint files)
111 files <- list.files("C:/Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project/input")
112
113 # write a function that takes the name of a file and returns the # of postive
114 # sentiment words, negative sentiment words, and the difference
115 GetSentiment <- function(file){
116 # get the file
117 fileName <- glue("input/", file, sep = "")
118 # get rid of any sneaky trailing spaces
119 fileName <- trimws(fileName)
120
121 df1 <- read.table(fileName, header=TRUE, fill = TRUE)
122
123 # read in the new file
124 fileText <- glue(read_file(fileName))
125 # remove any dollar signs (they're special characters in R)
126 fileText <- gsub("\\$", "", fileText)
127
128
129
130 # tokenize
131 tokens <- data_frame(text = fileText) %>% unnest_tokens(word, text)
132

```

```

133 # get the sentiment from the first text:
134 sentiment <- tokens %>%
135 inner_join(get_sentiments("bing")) %>% # pull out only sentiment words
136 count(sentiment) %>% # count the # of positive & negative words
137 spread(sentiment, n, fill = 0) %>% # made data wide rather than narrow
138 mutate(sentiment = positive - negative) %>% # # of positive words - # of negative words
139 mutate(file = file) %>% # add the name of our file
140 mutate(product = file) # add product
141
142
143
144 # return our sentiment dataframe
145 return(sentiment)
146 }
147
148 # file to put our output in
149 sentiments <- data_frame()
150
151 # get the sentiments for each file in our dataset
152 for(i in files){
153 sentiments <- rbind(sentiments, GetSentiment(i))
154 }
155
156 #plot overall sentiment for each product type
157 ggplot(sentiments, aes(x = file, y = sentiment, color = file)) +
158 geom_hline(yintercept = 0, color = "red") +
159 theme(axis.text.x = element_blank()) +
160 geom_boxplot() # draw a boxplot for each product
161
162
163 #Polarity by product
164 plot1 <- sentiments %>%
165 ggplot(aes(x = file, sentiment, fill = file)) +
166 geom_col() +
167 geom_hline(yintercept = 0, color = "red") +
168 theme(plot.title = element_text(size = 11)) +
169 xlab(NULL) + ylab(NULL) +
170 theme(axis.text.x = element_blank()) +
171 ggtitle("Polarity By Chart Level")
172 plot1
173
174
175 #plot nrc plot
176 undesirable_words <- c("dear", "sir", "madam")
177

```

```

178 #Store list of filenames (complaint files)
179 files <- list.files("C://Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project/input")
180
181 #glue path to file name
182 fileName <- glue("C://Users/sreel/OneDrive/XLRI/Term 4/BIDM/Project/input/", files[9], sep = "")
183
184 #Remove whitespace
185 fileName <- trimws(fileName)
186
187 fileText <- glue(read_file(fileName))
188
189 tokens <- data_frame(text = fileText) %>% unnest_tokens(word, text)
190 tokens %>%
191 filter(!inchar(word) < 3) %>% #remove words with less than 3 characters
192 filter(!word %in% undesirable_words) %>%
193 anti_join(stop_words)
194
195
196 savings_nrc <- tokens %>%
197 inner_join(get_sentiments("nrc"))
198
199 nrc_plot <- savings_nrc %>%
200 group_by(sentiment) %>%
201 summarise(word_count = n()) %>%
202 ungroup() %>%
203 mutate(sentiment = reorder(sentiment, word_count)) %>%
204 #Use `fill = -word_count` to make the larger bars darker
205 ggplot(aes(sentiment, word_count, fill = -word_count)) +
206 geom_col() +
207 guides(fill = FALSE) + #Turn off the legend
208 labs(x = NULL, y = "Word Count") +
209 scale_y_continuous(limits = c(0, 15000)) + #Hard code the axis limit
210 ggtitle("Vehicle loan or lease NRC Sentiment") +
211 coord_flip()
212
213 plot(nrc_plot)
214
215
216 #wordcloud
217 wordCloud_plot = function(file){
218 fileName <- glue("input/", file, sep = "")
219 # get rid of any sneaky trailing spaces
220 fileName <- trimws(fileName)
221

```

```

222 df1 <- read.table(fileName, header=TRUE, fill = TRUE)
223
224 # read in the new file
225 fileText <- glue(read_file(fileName))
226 # remove any dollar signs (they're special characters in R)
227 fileText <- gsub("\\$", "", fileText)
228
229
230
231 # tokenize
232 tokens <- data_frame(text = fileText) %>% unnest_tokens(word, text)
233
234 tokens %>%
235 inner_join(get_sentiments("bing")) %>%
236 count(word, sentiment, sort = TRUE) %>%
237 acast(word ~ sentiment, value.var = "n", fill = 0) %>%
238 comparison.cloud(colors = c("maroon", "darkgreen"),
239 max.words = 100)
240
241 }
242
243 wordCloud_plot(files[9])
244
245 #decision tree
246
247 library(dplyr)
248 dataset_main3 <- dataset_main
249 decision_data = select(dataset_main3, Product, Submitted.via, Company.response.to.consumer)
250 str(decision_data)
251
252 library(rpart)
253 library(rattle)

```

```

254 ac<-rpart(Company.response.to.consumer~.,data=decision_data,method="class",control =rpart.control(minsplit =1,minbucket=1, cp=0))
255 summary(ac)
256 plot(ac, uniform=TRUE, main="Classification tree for Company Response to consumers")
257 text(ac,pretty=0)
258
259 fancyRpartPlot(ac)
260
261
262 a<-as.data.frame(table(dataset_main$Product))
263 str(a)
264 pie(a$Freq,lbls = a$Var1)
265 a
266 library(ggplot2)
267 library(scales)
268 # Barplot
269 bp<- ggplot(a, aes(x="", y=a$Freq, fill=Var1))+
270 geom_bar(width = 1, stat = "identity")
271 bp
272 pie<- bp+coord_polar("y", start=0)+theme(axis.text.x = element_blank())+ xlab(NULL) + ylab(NULL)
273 pie
274
275 str(dataset_main)
276
277 #pie chart
278 p <- plot_ly(a, labels = ~Var1, values = ~Freq, type = 'pie',insidetextfont = list(color = '#FFFFFF'),marker = list(colors = colors,
279 line = list(color = '#FFFFFF', width = 1))) %>%
280 layout(title = 'Frequencies of Complaints by product type') %>%
281 layout(legend = list(x = 100, y = -10))
282 p
283
284
285 #plot complaint frequencies on the US map
286 library(maps)
287 library(mapttools)
288 library(sp)
289 library(plotly)
290 b<-as.data.frame(table(dataset_main$State))
291 str(b)
292 names(b)[1] <- 'state.abb'
293 str(b)
294 b$states <- tolower(state.name[match(b$state.abb, state.abb)])
295 str(b$states)
296 mapUSA <- map('state', fill = TRUE, plot = FALSE)
297 nms <- sapply(strsplit(mapUSA$names, ':'), function(x)x[1])

```

```

298 USApolygons <- map2SpatialPolygons(mapUSA, IDs = nms, CRS('+proj=longlat'))
299
300 idx <- match(unique(nms), b$states)
301 dat2 <- data.frame(value = b$Freq[idx], state = unique(nms))
302 row.names(dat2) <- unique(nms)
303
304 USAsp <- SpatialPolygonsDataFrame(USApolygons, data = dat2)
305 spplot(USAsp['value'])
306
307
308
309 #plot complaint frequencies by location
310 ca_data=subset(dataset_main, State == "CA")
311 ca_prod<-as.data.frame(table(ca_data$Product))
312 ca_prod_freq = ca_prod$Freq
313
314 tx_data=subset(dataset_main, State == "TX")
315 tx_prod<-as.data.frame(table(tx_data$Product))
316 tx_prod_freq = tx_prod$Freq
317
318 fl_data=subset(dataset_main, State == "FL")
319 fl_prod<-as.data.frame(table(fl_data$Product))
320 fl_prod_freq = fl_prod$Freq
321
322 products = ca_prod$Var1
323 data = data.frame(products,ca_prod_freq,tx_prod_freq)
324
325 p <- plot_ly(data, x = ~products, y = ~ca_prod_freq, type = 'bar', name = 'CA') %>%
326 add_trace(y = ~tx_prod_freq, name = 'TX') %>%
327 add_trace(y = ~fl_prod_freq, name = 'FL') %>%
328 layout(yaxis = list(title = 'Count'), barmode = 'stack')
329 p
330

```