

# Yelp: Exploratory Analysis

*Claire Jellison, Jacob Goldsmith, Ryan Kobler*

*11/20/2019*

## Overview:

The Yelp dataset includes over 5 million text reviews from businesses around the world. We aim to predict the number of stars a reviewer gives a business from the text of the review itself. To do so, we extract features of the text such as overall sentiment and word count to use as predictors.

Data source: [https://www.kaggle.com/yelp-dataset/yelp-dataset/version/6#yelp\\_review.csv](https://www.kaggle.com/yelp-dataset/yelp-dataset/version/6#yelp_review.csv)

## Packages:

```
library(jsonlite)
library(tidytext)
library(stringr)
library(wordcloud)

## Loading required package: RColorBrewer

library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##     filter, lag

## The following objects are masked from 'package:base':
##     intersect, setdiff, setequal, union

library(ggplot2)
library(GGally)

## Registered S3 method overwritten by 'GGally':
##     method from
##     +.gg    ggplot2

##
## Attaching package: 'GGally'

## The following object is masked from 'package:dplyr':
##     nasa
```

## Load in the data:

```
#path <- "/Users/ryankobler/Downloads/yelp_review.csv"
#yelp <- read.csv(path)

# Take sample of the data
yelp %>% sample_frac(0.01)
# Save sample for easy access
write.csv(df, "yelp-train.csv")

# Take sample of the data
yelp_sample <- yelp_review %>% sample_frac(0.01)
# Save sample for easy access
write.csv(yelp_sample, "yelp-train.csv")

# Load training data
yelp_train <- read.csv("yelp-train.csv")
yelp_sample <- yelp_train
```

## Clean

### Count total number of words

```
#count the number of words in the review
yelp_sample <- yelp_train
yelp_sample <- mutate(yelp_sample, numwords = str_count(yelp_sample$text, " "))
#univariate analysis of size
ggplot(data = yelp_sample, aes(x = numwords)) + geom_bar()
#analysis of size vs star rating with locally weighted polynomial
ggplot(data = yelp_sample, aes(x = numwords, y = stars)) + geom_jitter(size = 0.25) + geom_smooth()
```

### Define functions to extract features:

This method of feature analysis draws from the tidytext package and resource: <https://www.tidytextmining.com/sentiment.html>

```
# Note: this function requires the tidytext package & drops all
# words that do not convey sentiment
dropStopwords <- function(string){
  # Remove all punctuation except apostrophes & replace with " "
  noPunc <- gsub("[[:alnum:]][:space:]'", " ", string)
  # Split the larger string by space using strsplit()
  splitBySpace <- unlist(strsplit(noPunc, split = " "))
  # Remove missing chunks
  splitBySpace <- splitBySpace[splitBySpace != ""]
  todrop <- get_stopwords() # query dictionary of stopwords
  todrop <- todrop[[1]]

  # remove stop words and wrap as lowercase
  tolower(splitBySpace[!splitBySpace %in% todrop])
```

```

}

# Function below takes sentiment string in nrc and a sequence of
# pruned words associated with 1 review
nrcSentimentCount <- function(senti, text){
  sentiment <- nrc %>%
    filter(sentiment == senti) %>%
    select(word)

  # outputs a count of the number of "trues"
  sum(unlist(text) %in% unlist(sentiment))
}

# Remove the stop words and save in new column
#yelp_train$prunedtext <- lapply(yelp_train$text, FUN = dropStopwords)

get_sentiments("nrc") # we chose this one for now, can consider adding

## # A tibble: 13,901 x 2
##   word      sentiment
##   <chr>     <chr>
## 1 abacus    trust
## 2 abandon   fear
## 3 abandon   negative
## 4 abandon   sadness
## 5 abandoned anger
## 6 abandoned fear
## 7 abandoned negative
## 8 abandoned sadness
## 9 abandonment anger
## 10 abandonment fear
## # ... with 13,891 more rows

# other sentiment lexicons to increase our # of predictors
get_sentiments("afinn")

## # A tibble: 2,477 x 2
##   word      value
##   <chr>     <dbl>
## 1 abandon    -2
## 2 abandoned  -2
## 3 abandons   -2
## 4 abducted   -2
## 5 abduction  -2
## 6 abductions -2
## 7 abhor      -3
## 8 abhorred   -3
## 9 abhorrent  -3
## 10 abhors    -3
## # ... with 2,467 more rows

```

```
get_sentiments("loughran")

## # A tibble: 4,150 x 2
##   word      sentiment
##   <chr>     <chr>
## 1 abandon    negative
## 2 abandoned  negative
## 3 abandoning negative
## 4 abandonment negative
## 5 abandonments negative
## 6 abandons   negative
## 7 abdicated  negative
## 8 abdicates  negative
## 9 abdicating negative
## 10 abdication negative
## # ... with 4,140 more rows
```

### For knitting..

Because running the cleaning steps may take too long, we import the already clean data set now. The cleaning section below recounts our process of sampling from the larger data set and feature extraction.

```
yelp_train <- read.csv("DATA/withlanguage.csv")
yelp_train$prunedtext <- lapply(yelp_train$text, FUN = dropStopwords)

yelp_train <- yelp_train %>%
  filter(language == "english" | language == "scots")
```

### Data Description

The training data set, is a 1% random sample of the entire yelp\_review universe, made up of 52,616 observations and 19 variables. Each observation corresponds to a review from a random business. We use the NRC lexicon/dictionary to categorize all the non-stopwords in the reviews as anger, anticipation, disgust, fear, negative, positive, sadness, surprise, and trust. This was an arbitrary choice, for there are several other lexicons such as “loughran” and “afinn” both of which rank and categorize the sentiment of words differently.

The useful predictors include: - joy: number of joy-categorized words that appear in the body of the review (after removing stop words). And each of our variables named in this way follow the example given above.

We also include ratios for each, so `joyratio` is the total number of joy words divided by the total number of non-stop words to avoid privileging review length.

```
# Glimpse the data set to examine the predictors
glimpse(yelp_train)
```

```
## Observations: 51,945
## Variables: 33
## $ X                  <int> 1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, ...
## $ X1                 <int> 1, 2, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, ...
## $ review_id          <fct> 10c-gTpyEHIBeH6qG7iekw, 864AWE6eqipNE7SPVYSG...
## $ user_id            <fct> hyMa8iiLndcfqF3qJdN67w, cMEtAiW60I5wE_vLfTxo...
## $ business_id        <fct> grZEBaSzwWA3yJMwDR10Nw, LR1JNvpNrx2N4HNAYhHv...
```

```

## $ stars          <int> 4, 5, 2, 1, 2, 4, 3, 2, 5, 5, 4, 5, 3, 5, 5, ...
## $ date          <fct> 2011-02-20, 2014-08-11, 2016-12-14, 2017-03-...
## $ text          <fct> "I only came here for lunch so I can only sp...
## $ useful         <int> 3, 3, 0, 0, 0, 4, 2, 0, 0, 2, 2, 1, 0, 22, 0...
## $ funny          <int> 0, 2, 0, 0, 2, 5, 0, 0, 0, 0, 1, 0, 0, 14, 0...
## $ cool           <int> 2, 2, 0, 0, 0, 6, 0, 0, 0, 0, 1, 1, 0, 23, 0...
## $ joy            <int> 4, 6, 2, 3, 5, 6, 6, 8, 4, 11, 4, 9, 5, 3, 2...
## $ starsfactor    <int> 4, 5, 2, 1, 2, 4, 3, 2, 5, 5, 4, 5, 3, 5, 5, ...
## $ anger          <int> 0, 1, 2, 1, 5, 1, 5, 2, 0, 2, 0, 1, 1, 0, 0, ...
## $ nwords         <int> 75, 94, 50, 412, 231, 80, 192, 248, 18, 340, ...
## $ joyratio       <dbl> 0.053333333, 0.063829787, 0.040000000, 0.007...
## $ disgust        <int> 0, 1, 2, 1, 3, 2, 4, 1, 0, 3, 0, 0, 3, 0, 0, ...
## $ negative       <int> 0, 2, 3, 14, 6, 1, 5, 4, 1, 5, 2, 2, 4, 0, 2...
## $ positive       <int> 5, 12, 2, 18, 10, 8, 8, 12, 5, 20, 7, 10, 9, ...
## $ angerratio     <dbl> 0.000000000, 0.010638298, 0.040000000, 0.002...
## $ language        <fct> english, english, english, english, english, ...
## $ fear           <int> 0, 1, 1, 2, 3, 1, 1, 2, 1, 3, 1, 2, 4, 0, 2, ...
## $ surprise        <int> 2, 2, 0, 1, 2, 4, 4, 6, 1, 7, 1, 2, 3, 0, 2, ...
## $ trust           <int> 4, 4, 2, 13, 5, 6, 6, 8, 5, 15, 3, 6, 4, 4, ...
## $ anticipation   <int> 1, 5, 1, 17, 5, 5, 7, 5, 1, 14, 3, 3, 5, 3, ...
## $ fearratio      <dbl> 0.000000000, 0.010638298, 0.020000000, 0.004...
## $ positiveratio  <dbl> 0.066666667, 0.12765957, 0.040000000, 0.043689...
## $ negativeratio  <dbl> 0.000000000, 0.021276596, 0.060000000, 0.033...
## $ surpriseratio  <dbl> 0.026666667, 0.021276596, 0.000000000, 0.002...
## $ disgustratio   <dbl> 0.000000000, 0.010638298, 0.040000000, 0.002...
## $ trustratio     <dbl> 0.053333333, 0.042553191, 0.040000000, 0.031...
## $ anticipationratio <dbl> 0.013333333, 0.053191489, 0.020000000, 0.041...
## $ prunedtext      <list> [<"i", "came", "lunch", "i", "can", "speak"...

```

## Testing on tiny data set

```

yelp.small <- yelp_train[1:3,]
yelp.small$prunedtext <- lapply(yelp.small$text, FUN = dropStopwords)

yelp.small$prunedtext

# Matching sentiments from the NRC dictionary
nrc$sentiment <- as.factor(nrc$sentiment)

yelp.small %>%
  mutate(joy = nrcSentimentCount("fear", prunedtext))

# Grab all of the sentiments from for testing
sentiments <- levels(as.factor(nrc$sentiment))

# Generate new columns that count the number of times each sentiment word appears
yelp.small$joy <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "joy")
yelp.small$anger <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "anger")
yelp.small$fear <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "fear")
yelp.small$positive <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "positive")
yelp.small$negative <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "negative")
yelp.small$surprise <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "surprise")

```

```

yelp.small$disgust <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "disgust")
yelp.small$trust <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "trust")
yelp.small$anticipation <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "anticipation")
yelp.small$anger <- sapply(yelp.small$prunedtext, nrcSentimentCount, senti = "anger")

```

### Apply nrcSentimentCount function to the full training data set

Q: Is there an easier/cleaner way to create all of these columns?

```

# Generate new column called prunedtext
yelp_train$prunedtext <- lapply(yelp_train$text, FUN = dropStopwords)
# This column gives us a little trickiness if we try to export as .csv
class(yelp_train$prunedtext) # note that this column is list of lists

yelp_train$joy <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "joy")
yelp_train$anger <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "anger")
yelp_train$fear <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "fear")
yelp_train$positive <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "positive")
yelp_train$negative <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "negative")
yelp_train$surprise <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "surprise")
yelp_train$disgust <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "disgust")
yelp_train$trust <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "trust")
yelp_train$anticipation <- sapply(yelp_train$prunedtext, nrcSentimentCount, senti = "anticipation")

```

### Convert columns to numbers and normalize by length of review

Note that as.numeric is no longer explicitly necessary because we've used sapply instead of lapply to generate the sentiment counts. But it is not harmful.

```

# Add column that counts the number of total words
yelp_train$nwords <- str_count(yelp_train$text, " ")

# Generate the ratios
yelp_train$joyratio <- as.numeric(yelp_train$joy)/(yelp_train$nwords)
yelp_train$angerratio <- as.numeric(yelp_train$anger)/(yelp_train$nwords)
yelp_train$fearratio <- as.numeric(yelp_train$fear)/(yelp_train$nwords)
yelp_train$positiveratio <- as.numeric(yelp_train$positive)/(yelp_train$nwords)
yelp_train$negativeratio <- as.numeric(yelp_train$negative)/(yelp_train$nwords)
yelp_train$surpriseratio <- as.numeric(yelp_train$surprise)/(yelp_train$nwords)
yelp_train$disgustratio <- as.numeric(yelp_train$disgust)/(yelp_train$nwords)
yelp_train$strustratio <- as.numeric(yelp_train$trust)/(yelp_train$nwords)
yelp_train$anticipationratio <- as.numeric(yelp_train$anticipation)/(yelp_train$nwords)

```

### Missingness:

In this stage of the process, we found that some of our data was missing in that some reviews had no words detected. The problem came down to review language, so we use the `textcat` package to identify the language of the review and `dplyr` to filter out the non-English reviews.

This can be thought of as missingness in that our results only represent Yelp reviewers writing in English.

```

yelp_train %>%
  filter(nwords==0)

# This will take a while to run --> do overnight
# partition the data to see where the runtime issue is... the 3000 observation
# samples run in under a minute.
n <- nrow(yelp_train)
yelp_train$language <- NA
yelp_train[1:1000,]$language <- textcat(yelp_train[1:1000,]$text)
yelp_train[1000:3000,]$language <- textcat(yelp_train[1000:3000,]$text)
yelp_train[3000:6000,]$language <- textcat(yelp_train[3000:6000,]$text)
yelp_train[6001:9000,]$language <- textcat(yelp_train[6001:9000,]$text)
yelp_train[9000:40000,]$language <- textcat(yelp_train[9000:40000,]$text)
yelp_train[40000:n,]$language <- textcat(yelp_train[40000:n,]$text)

# save languages in another csv
langs <- yelp_train %>%
  select(X1, review_id, user_id, business_id, stars, language)

# writing only id vars and languages
write.csv(langs, "DATA/withlanguage.csv")

# Remove non-English reviews
yelp_train_en <- yelp_train %>%
  filter(language == "english" | language == "scots")

# Dropped 671 observations
nrow(yelp_train)
nrow(yelp_train_en)

```

### Write yelp\_train data frame to a csv

Save in /DATA so we don't have to run everything again. We remove the prunedtext column because we're guessing that the csv filetype cannot handle list type entries.

```

write.csv(yelp_train %>% select(-prunedtext), "DATA/withlanguage.csv")
write.csv(yelp_train[, -12], "yelp-train3.csv")

```

### Generate word clouds for each rating level

```

# Save individual words from the X-star reviews as the vector wordsX
words4 <- yelp_train %>%
  filter(stars == 4) %>%
  pull(prunedtext)

words5 <- yelp_train %>%
  filter(stars == 5) %>%
  pull(prunedtext)

words1 <- yelp_train %>%

```

```

filter(stars == 1) %>%
pull(prunedtext)

# we unlist() the above vectors to get big pile o'
# words across all reviews since wordcloud() takes character vectors as input,
# not lists
wordcloud(unlist(words5), min.freq = 10, max.words = 30)

```

## Loading required namespace: tm

## Warning in tm\_map.SimpleCorpus(corpus, tm::removePunctuation):  
## transformation drops documents

## Warning in tm\_map.SimpleCorpus(corpus, function(x) tm::removeWords(x,  
## tm::stopwords())): transformation drops documents



```
wordcloud(unlist(words1), min.freq = 10, max.words = 30)
```

```

## Warning in tm_map.SimpleCorpus(corpus, tm::removePunctuation):
## transformation drops documents

## Warning in tm_map.SimpleCorpus(corpus, tm::removePunctuation):
## transformation drops documents

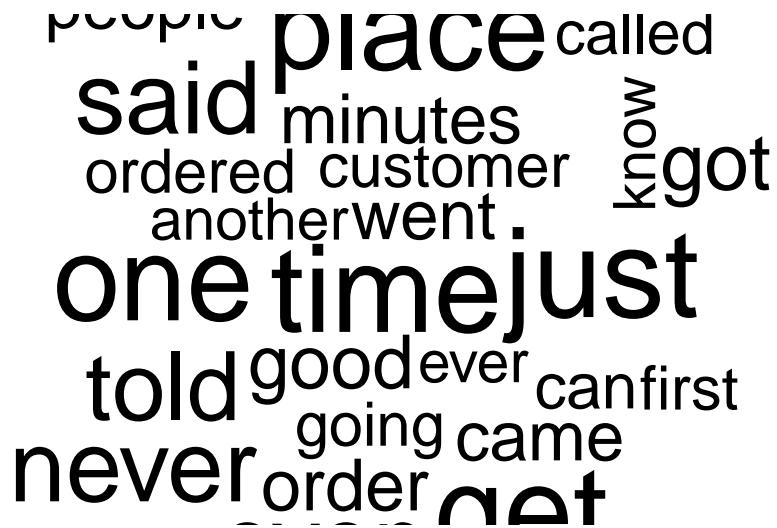
```

```

## Warning in wordcloud(unlist(words1), min.freq = 10, max.words = 30): back
## could not be fit on page. It will not be plotted.

## Warning in wordcloud(unlist(words1), min.freq = 10, max.words = 30):
## service could not be fit on page. It will not be plotted.

```



Idea: could think about coloring the words by sentiment.

### Histograms of Sentiments

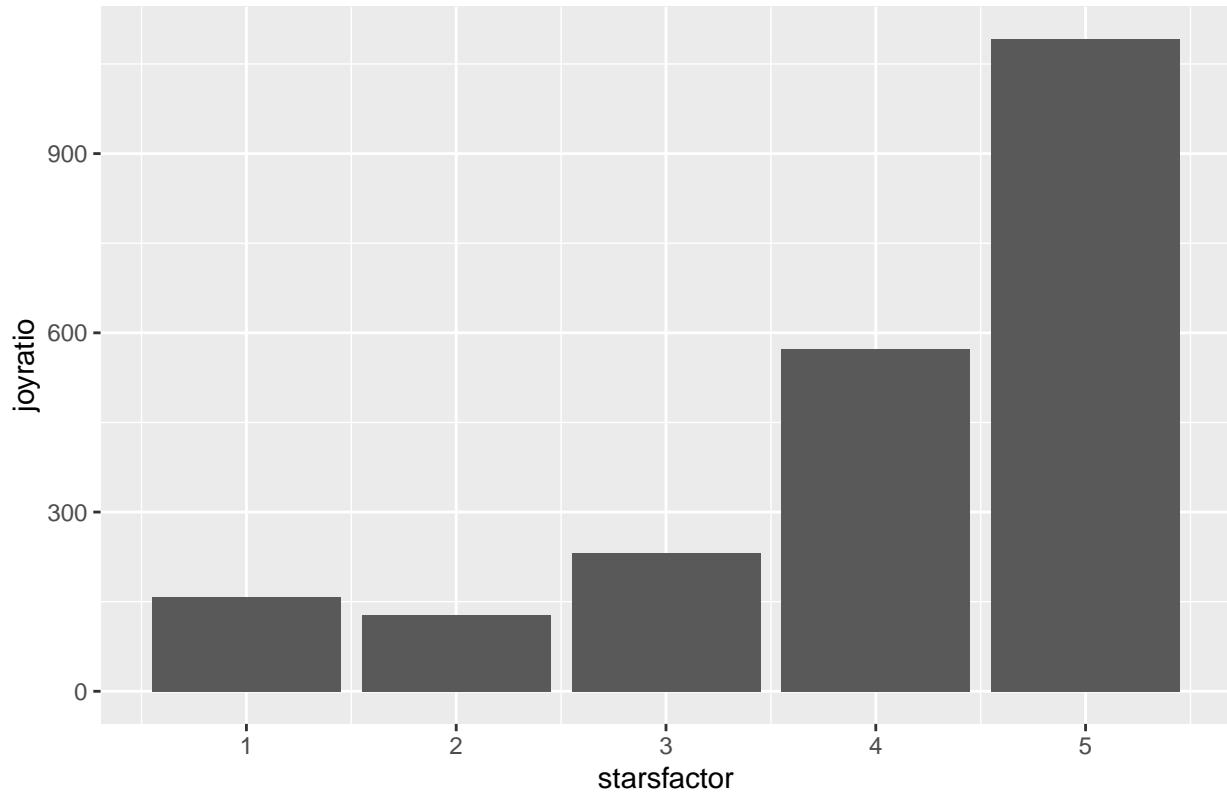
Bar chart of ratio of ‘joy’ to total number of words.

```

# Joy
ggplot(yelp_train, aes(starsfactor, joyratio)) +
  geom_bar(stat = "identity") +
  ggtitle("Distribution of joy words")

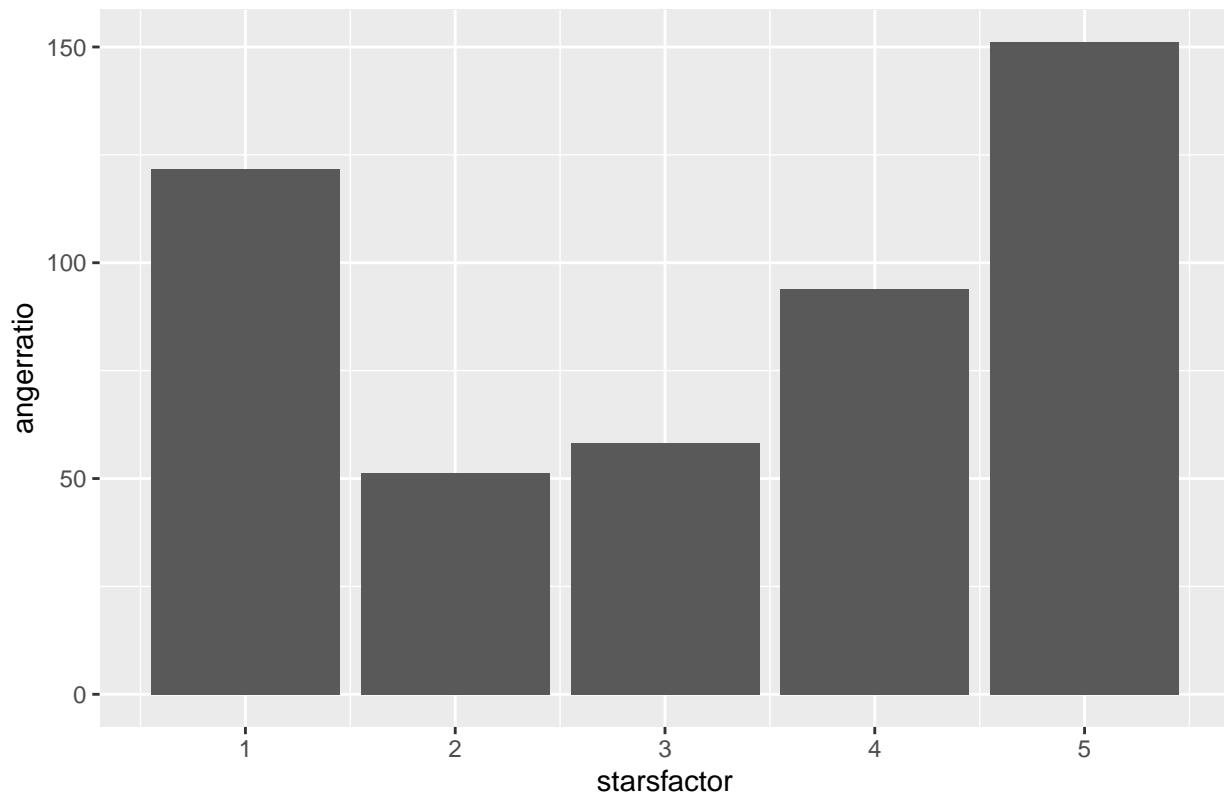
```

## Distribution of joy words



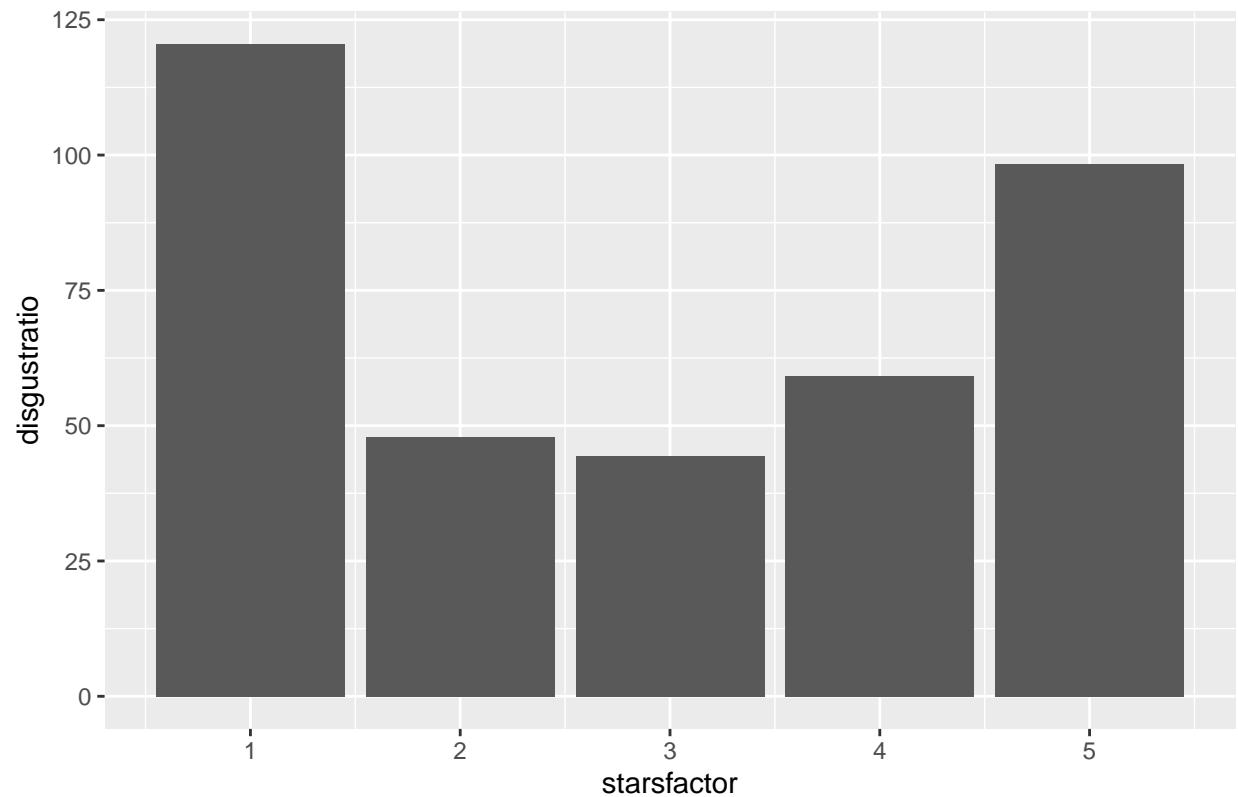
```
# Anger
ggplot(yelp_train, aes(starsfactor, angerratio)) +
  geom_bar(stat = "identity") +
  ggtitle("Distribution of anger words")
```

Distribution of anger words



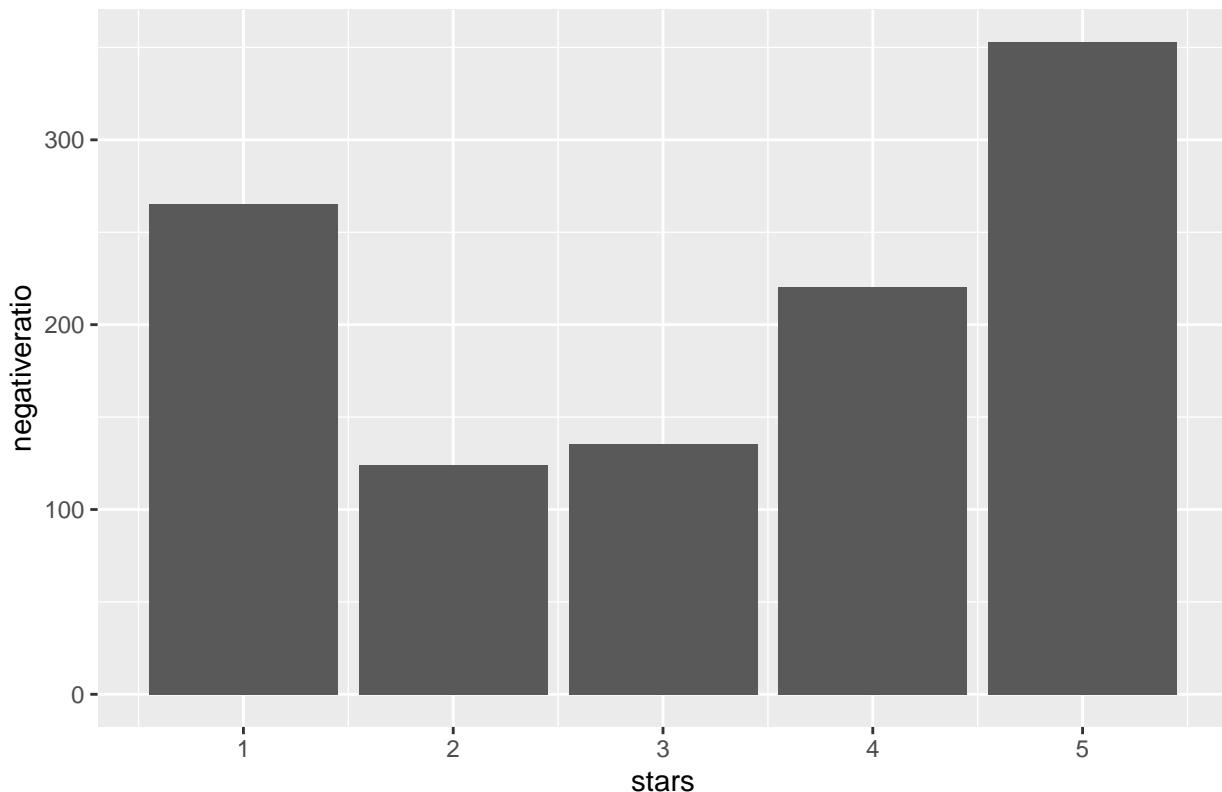
```
# Disgust
ggplot(yelp_train, aes(starsfactor, disgustratio)) +
  geom_bar(stat = "identity") +
  ggtitle("Distribution of disgust words")
```

### Distribution of disgust words



```
# Negativity
ggplot(yelp_train, aes(stars, negativeratio)) +
  geom_bar(stat = "identity") +
  ggtitle("Distribution of negative words")
```

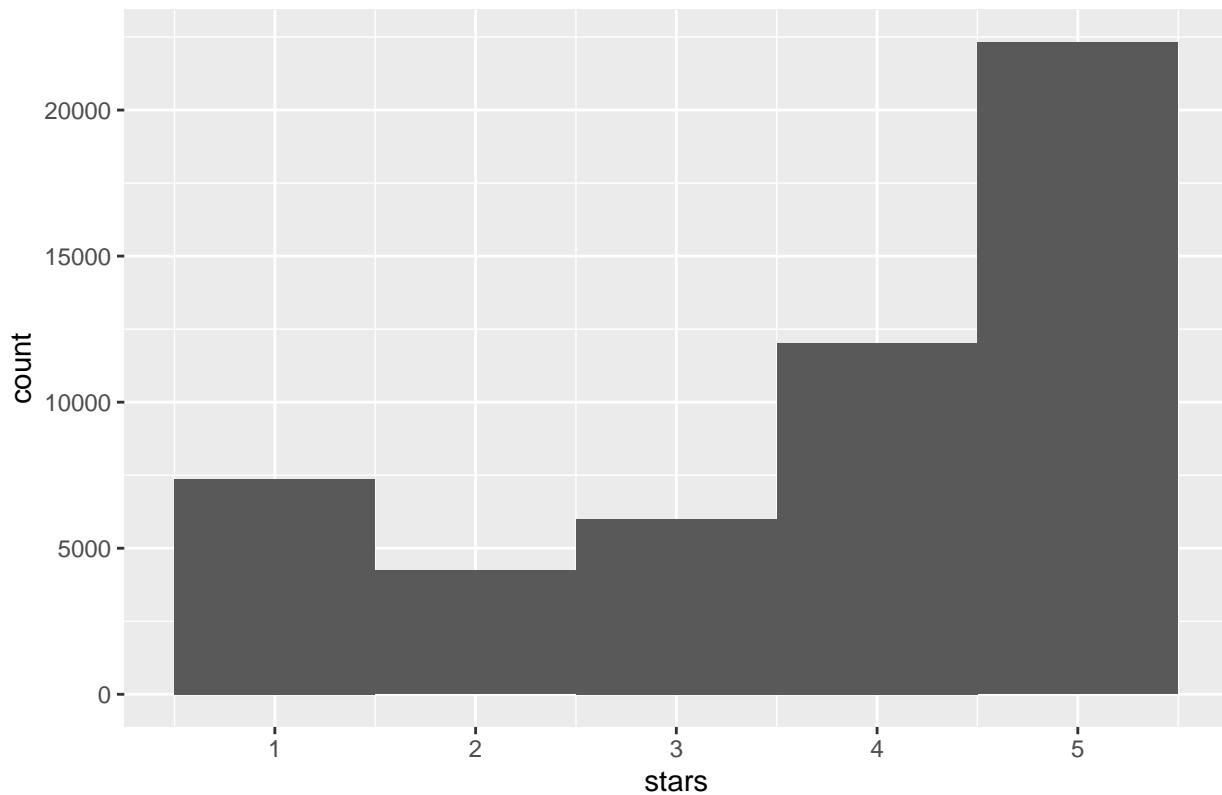
## Distribution of negative words



## Univariate analysis of the response

```
mean(yelp_train$stars)  
  
## [1] 3.725575  
  
var(yelp_train$stars)  
  
## [1] 2.072406  
  
ggplot(yelp_train, aes(x=stars)) +  
  geom_histogram(binwidth=1) +  
  ggtitle("Bar Graph of Stars")
```

## Bar Graph of Stars

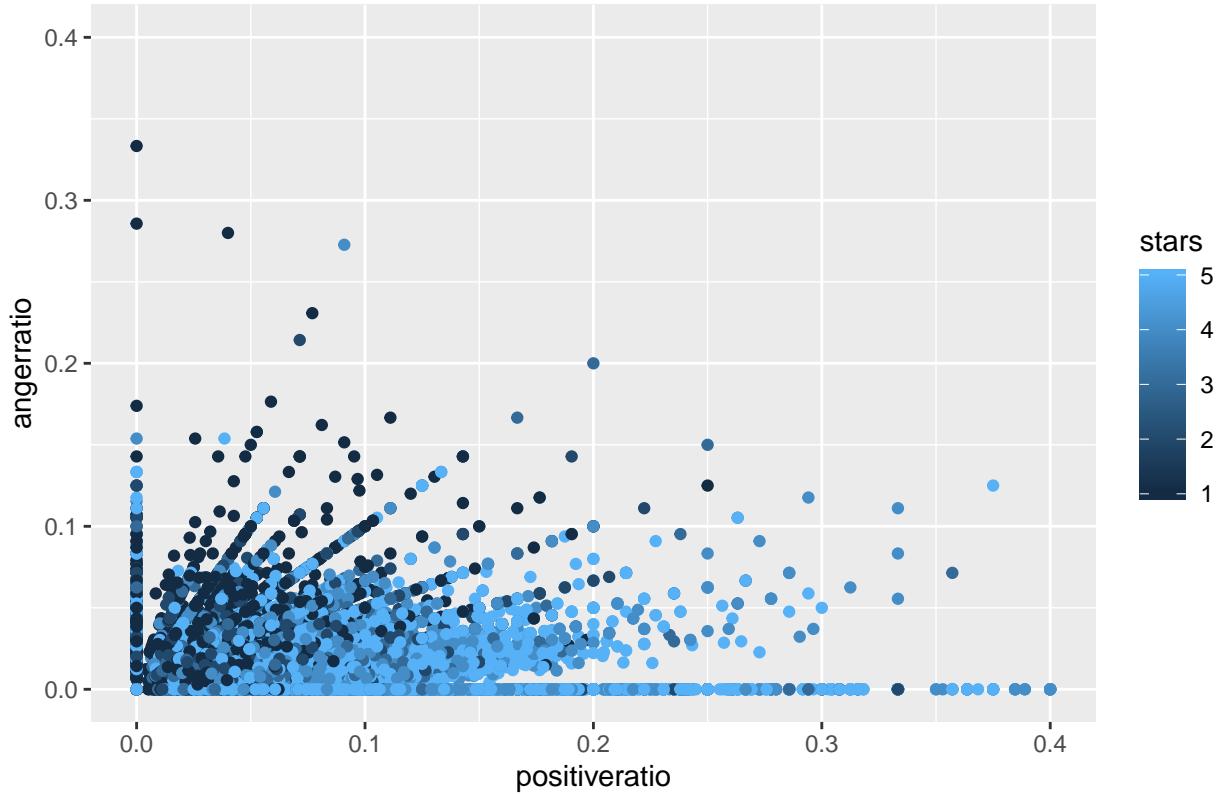


The mean number of stars in our sample is around 3.726 and the variance is around 2.067. We can see that the distribution of stars is somewhat oddly shaped with lots of five and four star review and a fair number of 1 star reviews.

```
ggplot(yelp_train, aes(x=positiveratio, y = angerratio)) +  
  geom_point(aes(color = stars)) +  
  ggtitle("Scatter of Stars") + xlim(x = 0, .4) + ylim(y = 0, .4)
```

```
## Warning: Removed 29 rows containing missing values (geom_point).
```

## Scatter of Stars



As expected, it appears from the graph that words labeled as “anger” tend to have a smaller number of stars whereas words that are labeled as positive tend to have a higher number of stars. We can also see that positive words are more frequently used in higher proportion than angry words.

Even though each of the sentiments are normalized by the number of words total, there may simply be more positive words in the NRC dictionary, increasing the likelihood that any one review finds positive words that match. In fitting a model, we will want to be cognizant of this.

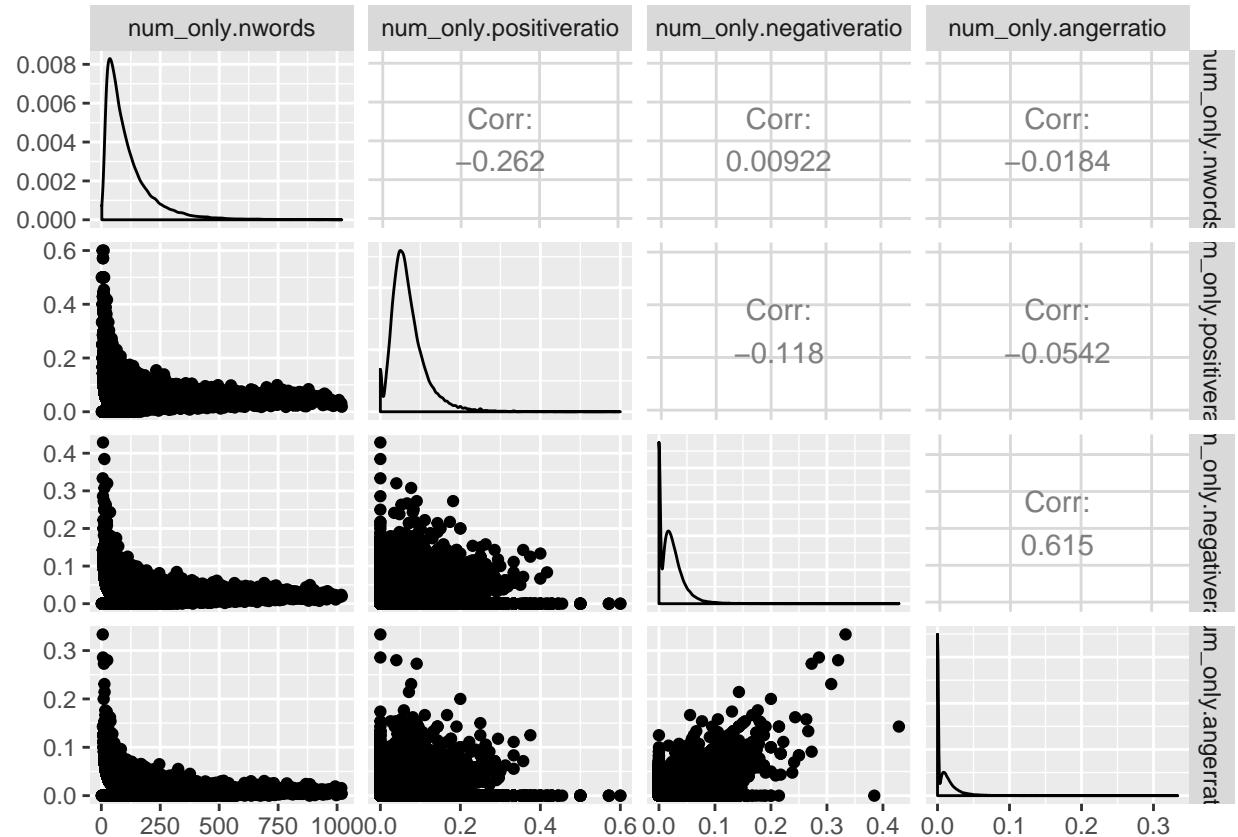
## Bivariate/Trivariate Graphs

```

bool <- sapply(yelp_train, is.numeric)
num_only <- yelp_train[,bool]
#scatterplot matrix of all of the variables
#ggpairs(num_only, labels = colnames(num_only)) # this is not very useful
#so I create a data.frame with only 4 essential variables
selected <- data.frame(num_only$nwords, num_only$positiveratio, num_only$negativeratio, num_only$angerratio)
#do a corrrgram on those vars
ggpairs(selected, labels = colnames(selected))

## Warning in warn_if_args_exist(list(...)): Extra arguments: 'labels' are
## being ignored. If these are meant to be aesthetics, submit them using the
## 'mapping' variable within ggpairs with ggplot2::aes or ggplot2::aes_string.

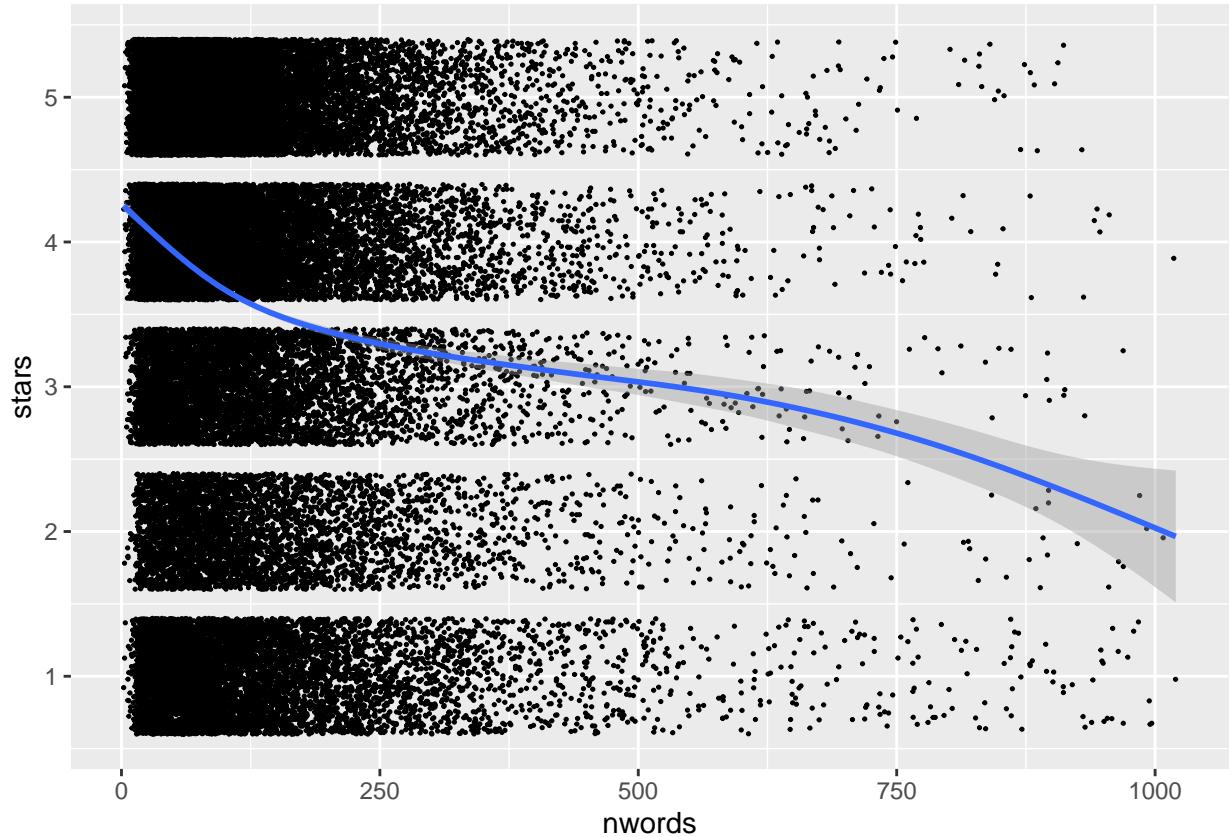
```



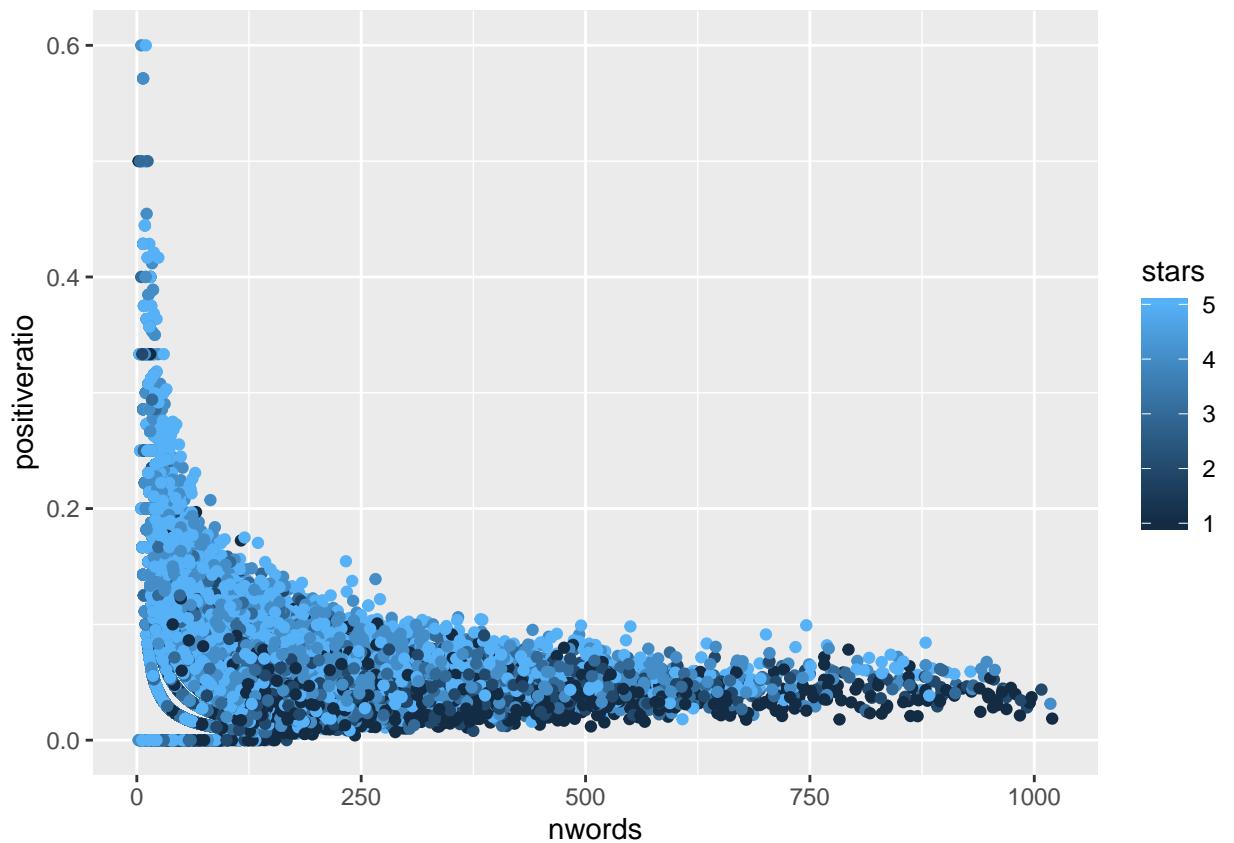
#selected scatterplots of important variables

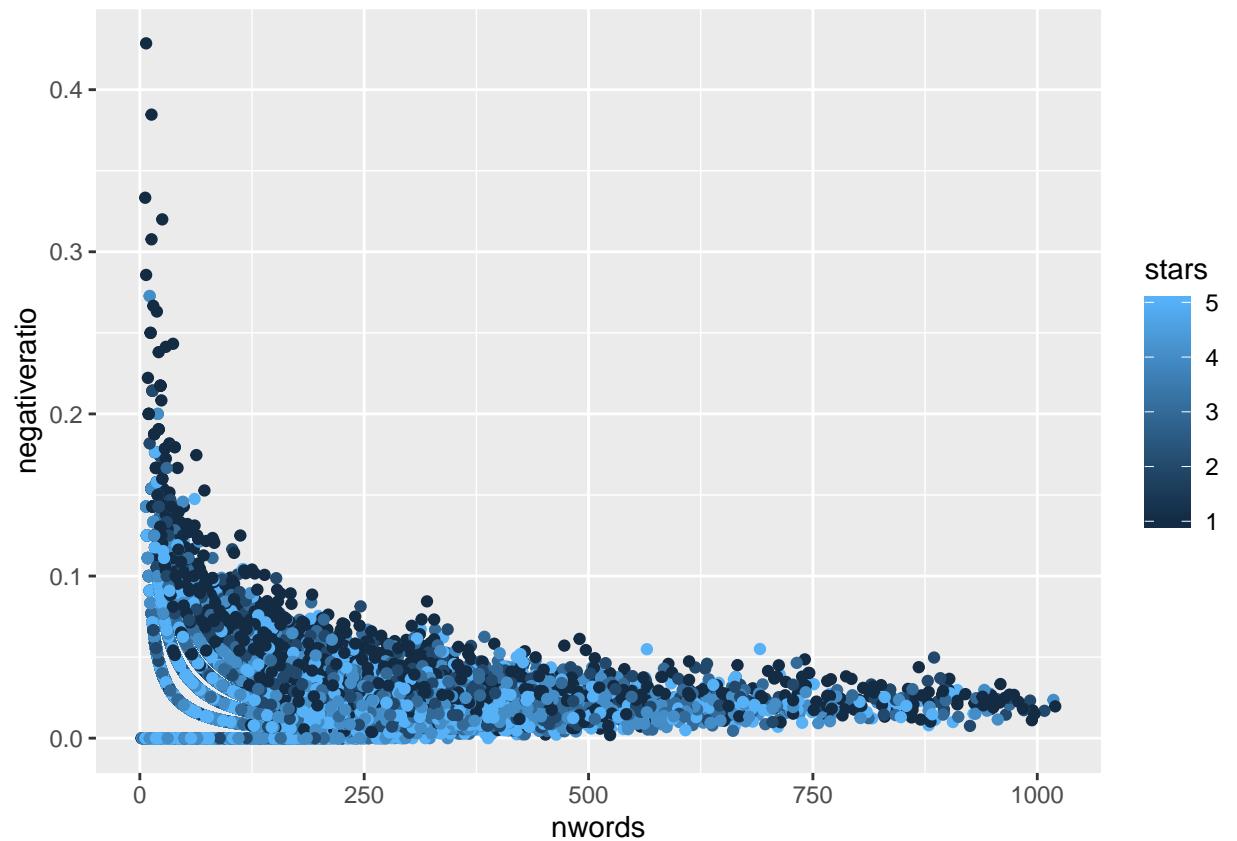
```
ggplot(data = num_only, aes(x = nwords, y = stars)) + geom_jitter(size = 0.25) + geom_smooth() #stars
```

```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```

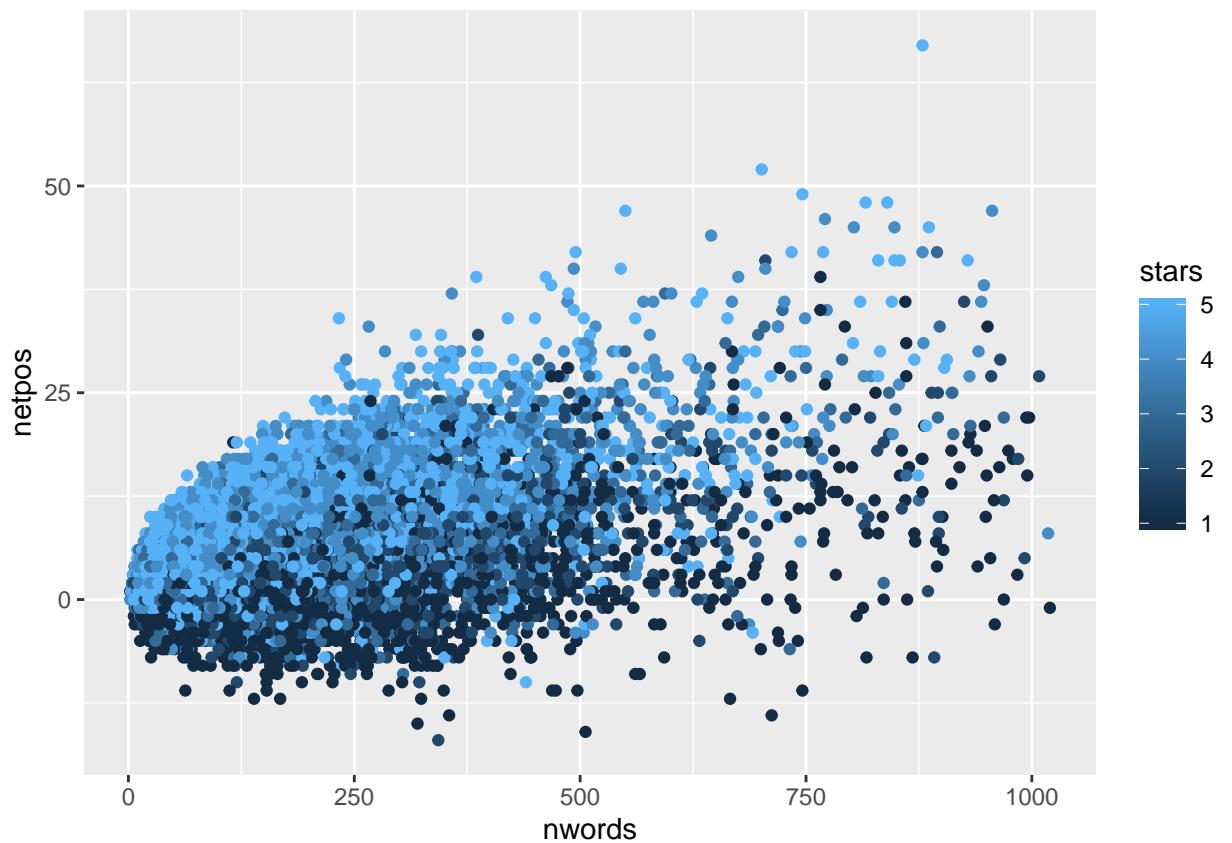


```
ggplot(data = num_only, aes(x = nwords, y = positiveratio)) + geom_point(aes(color = stars))
```

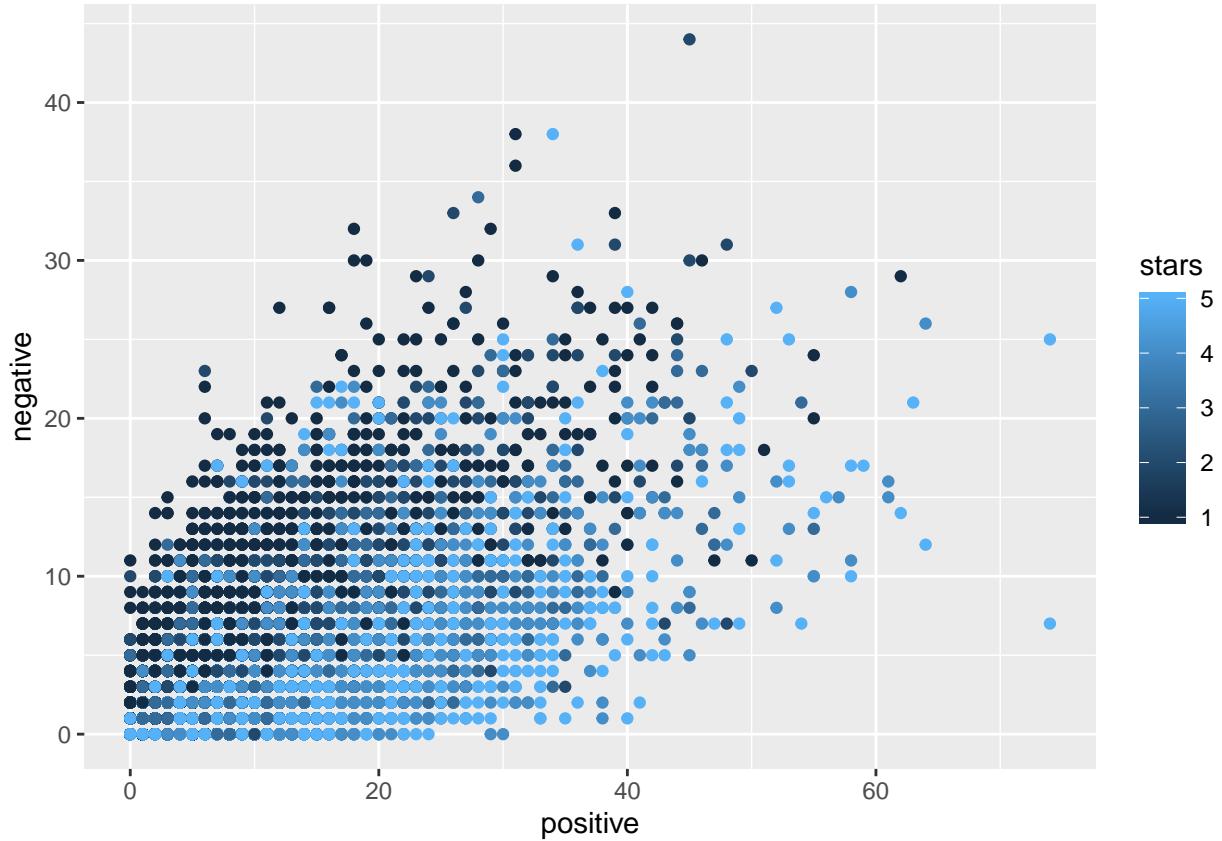




```
#adding a column that is positive - negative
num_only$netpos <- num_only$positive - num_only$negative
ggplot(data = num_only, aes(x = nwords, y = netpos)) + geom_point(aes(color = stars))
```



```
ggplot(data = num_only, aes(x = positive, y = negative)) + geom_point(aes(color = stars))
```



```
cor(num_only, use = "complete.obs")
```

	X	X1	stars	useful
## X	1.000000e+00	1.000000e+00	-0.000663017	0.0005636345
## X1	1.000000e+00	1.000000e+00	-0.000663017	0.0005636345
## stars	-6.630170e-04	-6.630170e-04	1.000000000	-0.1044024865
## useful	5.636345e-04	5.636345e-04	-0.104402486	1.000000000
## funny	-3.716674e-03	-3.716674e-03	-0.050009023	0.5215216579
## cool	-1.752300e-03	-1.752300e-03	0.048583885	0.6838093548
## joy	-4.047080e-03	-4.047080e-03	0.056999857	0.1738751646
## starsfactor	-6.630170e-04	-6.630170e-04	1.000000000	-0.1044024865
## anger	-6.375435e-06	-6.375435e-06	-0.288100319	0.2093829374
## nwords	-6.331144e-03	-6.331144e-03	-0.192387823	0.2807878195
## joyratio	2.390498e-03	2.390498e-03	0.260511359	-0.1096721967
## disgust	-6.186219e-04	-6.186219e-04	-0.354046547	0.1773691831
## negative	-6.684311e-03	-6.684311e-03	-0.332064351	0.2431809639
## positive	-4.404313e-03	-4.404313e-03	-0.015199413	0.2316128361
## angerratio	4.196042e-03	4.196042e-03	-0.224083146	0.0266812909
## fear	-3.304916e-03	-3.304916e-03	-0.239043777	0.2072699366
## surprise	-2.729834e-03	-2.729834e-03	-0.062910699	0.1791594478
## trust	-2.721242e-03	-2.721242e-03	-0.045166851	0.2165931354
## anticipation	-6.025198e-03	-6.025198e-03	-0.081234108	0.2163158860
## fearratio	4.186770e-03	4.186770e-03	-0.172136931	0.0260869301
## positiveratio	3.085654e-03	3.085654e-03	0.277181375	-0.0992574470
## negativeratio	-5.272395e-04	-5.272395e-04	-0.314749000	0.0356770342

	4.143237e-04	4.143237e-04	0.069314959	-0.0430573111
## surpriseratio	5.183515e-03	5.183515e-03	-0.300331361	0.0194803605
## disgustratio	4.457505e-03	4.457505e-03	0.188683148	-0.0814291991
## trustratio	7.910009e-04	7.910009e-04	0.116064101	-0.0547205894
## anticipationratio	-1.110814e-03	-1.110814e-03	0.198146470	0.1292277594
## netpos				
## funny		cool	joy	starsfactor
## X	-0.003716674	-0.001752300	-0.00404708	-0.000663017
## X1	-0.003716674	-0.001752300	-0.00404708	-0.000663017
## stars	-0.050009023	0.048583885	0.05699986	1.000000000
## useful	0.521521658	0.683809355	0.17387516	-0.104402486
## funny	1.000000000	0.606993947	0.11348697	-0.050009023
## cool	0.606993947	1.000000000	0.19365965	0.048583885
## joy	0.113486971	0.193659646	1.000000000	0.056999857
## starsfactor	-0.050009023	0.048583885	0.05699986	1.000000000
## anger	0.126984282	0.124641701	0.38720432	-0.288100319
## nwords	0.157824904	0.208945862	0.68630151	-0.192387823
## joyratio	-0.050699007	-0.046890914	0.24728050	0.260511359
## disgust	0.115521870	0.095580087	0.33628983	-0.354046547
## negative	0.151245967	0.148691836	0.46786665	-0.332064351
## positive	0.132727986	0.209920324	0.87968837	-0.015199413
## angerratio	0.021031992	-0.005362692	-0.04096215	-0.224083146
## fear	0.122481628	0.122638322	0.38115206	-0.239043777
## surprise	0.115136653	0.164932636	0.73250737	-0.062910699
## trust	0.119200641	0.178642922	0.84942773	-0.045166851
## anticipation	0.121471643	0.184235860	0.77552365	-0.081234108
## fearratio	0.018304406	-0.006568931	-0.04743176	-0.172136931
## positiveratio	-0.052302345	-0.045581079	0.16074823	0.277181375
## negativeratio	0.032654392	-0.004574759	-0.05883654	-0.314749000
## surpriseratio	-0.014497628	-0.015946771	0.16688047	0.069314959
## disgustratio	0.021414680	-0.014757186	-0.05044931	-0.300331361
## trustratio	-0.046630858	-0.050345291	0.15227675	0.188683148
## anticipationratio	-0.029213614	-0.023075049	0.15150822	0.116064101
## netpos	0.066283974	0.164000785	0.78873271	0.198146470
##				
## anger		nwords	joyratio	disgust
## X	-6.375435e-06	-0.006331144	0.002390498	-0.0006186219
## X1	-6.375435e-06	-0.006331144	0.002390498	-0.0006186219
## stars	-2.881003e-01	-0.192387823	0.260511359	-0.3540465475
## useful	2.093829e-01	0.280787819	-0.109672197	0.1773691831
## funny	1.269843e-01	0.157824904	-0.050699007	0.1155218702
## cool	1.246417e-01	0.208945862	-0.046890914	0.0955800873
## joy	3.872043e-01	0.686301515	0.247280501	0.3362898309
## starsfactor	-2.881003e-01	-0.192387823	0.260511359	-0.3540465475
## anger	1.000000e+00	0.608426860	-0.180438400	0.6987039685
## nwords	6.084269e-01	1.000000000	-0.260208519	0.5540852919
## joyratio	-1.804384e-01	-0.260208519	1.000000000	-0.1787461688
## disgust	6.987040e-01	0.554085292	-0.178746169	1.0000000000
## negative	7.764983e-01	0.770075167	-0.239146316	0.7320104769
## positive	4.804277e-01	0.837669874	0.023072280	0.4202136214
## angerratio	5.282159e-01	-0.018448589	-0.039961142	0.3230421061
## fear	6.950502e-01	0.615017668	-0.187336404	0.6112795179
## surprise	4.526303e-01	0.647675208	0.052238539	0.3874181561
## trust	4.744537e-01	0.778497147	0.065547612	0.4185024640
## anticipation	4.754848e-01	0.778319586	-0.016417874	0.4067394821
## fearratio	2.807417e-01	-0.013082221	-0.047928950	0.2477467185

```

## positiveratio -1.914642e-01 -0.262213602 0.834365959 -0.1913223761
## negativeratio 3.413719e-01 0.009220264 -0.102947303 0.3472330639
## surpriseratio -1.255243e-02 -0.098018698 0.507539340 -0.0340846996
## disgustratio 3.020755e-01 -0.014633329 -0.062744477 0.5589056959
## trustratio -1.404445e-01 -0.223569761 0.779334808 -0.1379918818
## anticipationratio -8.636859e-02 -0.139733940 0.555250547 -0.1007189151
## netpos 9.020762e-02 0.538908018 0.185029207 0.0443657950
## negative positive angerratio fear
## X -0.006684311 -0.004404313 0.004196042 -0.003304916
## X1 -0.006684311 -0.004404313 0.004196042 -0.003304916
## stars -0.332064351 -0.015199413 -0.224083146 -0.239043777
## useful 0.243180964 0.231612836 0.026681291 0.207269937
## funny 0.151245967 0.132727986 0.021031992 0.122481628
## cool 0.148691836 0.209920324 -0.005362692 0.122638322
## joy 0.467866655 0.879688372 -0.040962151 0.381152064
## starsfactor -0.332064351 -0.015199413 -0.224083146 -0.239043777
## anger 0.776498250 0.480427670 0.528215929 0.695050215
## nwords 0.770075167 0.837669874 -0.018448589 0.615017668
## joyratio -0.239146316 0.023072280 -0.039961142 -0.187336404
## disgust 0.732010477 0.420213621 0.323042106 0.611279518
## negative 1.000000000 0.599670783 0.243956978 0.722040899
## positive 0.599670783 1.000000000 -0.043426018 0.491502130
## angerratio 0.243956978 -0.043426018 1.000000000 0.264813040
## fear 0.722040899 0.491502130 0.264813040 1.000000000
## surprise 0.517411124 0.708845709 0.038618675 0.438036218
## trust 0.574535097 0.887342656 -0.020364039 0.482784652
## anticipation 0.597270959 0.810160212 -0.016812499 0.477981882
## fearratio 0.200732638 -0.039243817 0.543668853 0.529044916
## positiveratio -0.247064651 0.102355780 -0.054199737 -0.189901751
## negativeratio 0.442559824 -0.045119909 0.614625503 0.267687235
## surpriseratio -0.057802004 0.051029161 0.081929173 -0.025954712
## disgustratio 0.234202175 -0.049516708 0.623835501 0.218603833
## trustratio -0.195257886 0.031448861 -0.012029461 -0.140399356
## anticipationratio -0.105946313 0.039766468 -0.004861244 -0.086524900
## netpos 0.092481922 0.852275947 -0.213499087 0.139583923
## surprise trust anticipation fearratio
## X -0.002729834 -0.002721242 -0.0060251979 4.186770e-03
## X1 -0.002729834 -0.002721242 -0.0060251979 4.186770e-03
## stars -0.062910699 -0.045166851 -0.0812341084 -1.721369e-01
## useful 0.179159448 0.216593135 0.2163158860 2.608693e-02
## funny 0.115136653 0.119200641 0.1214716433 1.830441e-02
## cool 0.164932636 0.178642922 0.1842358602 -6.568931e-03
## joy 0.732507375 0.849427729 0.7755236527 -4.743176e-02
## starsfactor -0.062910699 -0.045166851 -0.0812341084 -1.721369e-01
## anger 0.452630348 0.474453675 0.4754847770 2.807417e-01
## nwords 0.647675208 0.778497147 0.7783195859 -1.308222e-02
## joyratio 0.052238539 0.065547612 -0.0164178741 -4.792895e-02
## disgust 0.387418156 0.418502464 0.4067394821 2.477467e-01
## negative 0.517411124 0.574535097 0.5972709586 2.007326e-01
## positive 0.708845709 0.887342656 0.8101602119 -3.924382e-02
## angerratio 0.038618675 -0.020364039 -0.0168124990 5.436689e-01
## fear 0.438036218 0.482784652 0.4779818818 5.290449e-01
## surprise 1.000000000 0.699393184 0.7483493246 2.891678e-02
## trust 0.699393184 1.000000000 0.7927254374 -1.690089e-02

```

## anticipation	0.748349325	0.792725437	1.0000000000	-1.303523e-02
## fearratio	0.028916778	-0.016900891	-0.0130352295	1.000000e+00
## positiveratio	0.007618102	0.067026806	-0.0375465397	-5.402303e-02
## negativeratio	0.028831849	-0.027893342	0.0064396535	4.852689e-01
## surpriseratio	0.437635688	0.098863957	0.1347740593	6.741620e-02
## disgustratio	0.010083409	-0.026454452	-0.0319490281	4.690853e-01
## trustratio	0.045119788	0.208737408	-0.0004659495	-1.592049e-02
## anticipationratio	0.172182264	0.092207153	0.3002244921	6.655761e-06
## netpos	0.543774743	0.728531313	0.6176347437	-1.800412e-01
##				
	positiveratio	negativeratio	surpriseratio	disgustratio
## X	0.003085654	-0.0005272395	0.0004143237	0.005183515
## X1	0.003085654	-0.0005272395	0.0004143237	0.005183515
## stars	0.277181375	-0.3147490000	0.0693149592	-0.300331361
## useful	-0.099257447	0.0356770342	-0.0430573111	0.019480360
## funny	-0.052302345	0.0326543922	-0.0144976275	0.021414680
## cool	-0.045581079	-0.0045747589	-0.0159467710	-0.014757186
## joy	0.160748229	-0.0588365355	0.1668804741	-0.050449308
## starsfactor	0.277181375	-0.3147490000	0.0693149592	-0.300331361
## anger	-0.191464177	0.3413719387	-0.0125524257	0.302075506
## nwords	-0.262213602	0.0092202640	-0.0980186982	-0.014633329
## joyratio	0.834365959	-0.1029473026	0.5075393400	-0.062744477
## disgust	-0.191322376	0.3472330639	-0.0340846996	0.558905696
## negative	-0.247064651	0.4425598243	-0.0578020035	0.234202175
## positive	0.102355780	-0.0451199086	0.0510291614	-0.049516708
## angerratio	-0.054199737	0.6146255032	0.0819291729	0.623835501
## fear	-0.189901751	0.2676872351	-0.0259547122	0.218603833
## surprise	0.007618102	0.0288318490	0.4376356883	0.010083409
## trust	0.067026806	-0.0278933416	0.0988639569	-0.026454452
## anticipation	-0.037546540	0.0064396535	0.1347740593	-0.031949028
## fearratio	-0.054023032	0.4852689042	0.0674162007	0.469085327
## positiveratio	1.0000000000	-0.1183609123	0.4204233483	-0.077101549
## negativeratio	-0.118360912	1.0000000000	0.0356875429	0.594496670
## surpriseratio	0.420423348	0.0356875429	1.0000000000	0.025420820
## disgustratio	-0.077101549	0.5944966701	0.0254208198	1.0000000000
## trustratio	0.763194120	-0.0760231149	0.4790273475	-0.030842120
## anticipationratio	0.497394703	-0.0011029457	0.5699456408	-0.040770347
## netpos	0.288854353	-0.3454262369	0.1012765809	-0.214701121
##				
	trustratio	anticipationratio	netpos	
## X	0.0044575047	7.910009e-04	-0.001110814	
## X1	0.0044575047	7.910009e-04	-0.001110814	
## stars	0.1886831484	1.160641e-01	0.198146470	
## useful	-0.0814291991	-5.472059e-02	0.129227759	
## funny	-0.0466308575	-2.921361e-02	0.066283974	
## cool	-0.0503452911	-2.307505e-02	0.164000785	
## joy	0.1522767527	1.515082e-01	0.788732711	
## starsfactor	0.1886831484	1.160641e-01	0.198146470	
## anger	-0.1404445099	-8.636859e-02	0.090207621	
## nwords	-0.2235697605	-1.397339e-01	0.538908018	
## joyratio	0.7793348081	5.552505e-01	0.185029207	
## disgust	-0.1379918818	-1.007189e-01	0.044365795	
## negative	-0.1952578856	-1.059463e-01	0.092481922	
## positive	0.0314488607	3.976647e-02	0.852275947	
## angerratio	-0.0120294612	-4.861244e-03	-0.213499087	
## fear	-0.1403993562	-8.652490e-02	0.139583923	

```

## surprise          0.0451197881    1.721823e-01  0.543774743
## trust            0.2087374078    9.220715e-02  0.728531313
## anticipation    -0.0004659495   3.002245e-01  0.617634744
## fearratio        -0.0159204856   6.655761e-06  -0.180041154
## positiveratio    0.7631941199    4.973947e-01  0.288854353
## negativeratio    -0.0760231149   -1.102946e-03 -0.345426237
## surpriseratio    0.4790273475    5.699456e-01  0.101276581
## disgustratio     -0.0308421205   -4.077035e-02 -0.214701121
## trustratio       1.0000000000    5.353805e-01  0.166763562
## anticipationratio 0.5353804504    1.000000e+00  0.118733069
## netpos           0.1667635617    1.187331e-01  1.000000000

```

*#on the smaller set*

```
cor(selected, use = "complete.obs")
```

```

##                               num_only.nwords num_only.positiveratio
## num_only.nwords             1.0000000000      -0.26221360
## num_only.positiveratio     -0.262213602       1.000000000
## num_only.negativeratio     0.009220264       -0.11836091
## num_only.angerratio        -0.018448589      -0.05419974
##                               num_only.negativeratio num_only.angerratio
## num_only.nwords              0.009220264      -0.01844859
## num_only.positiveratio      -0.118360912      -0.05419974
## num_only.negativeratio      1.0000000000      0.61462550
## num_only.angerratio         0.614625503      1.000000000

```

```
```

```