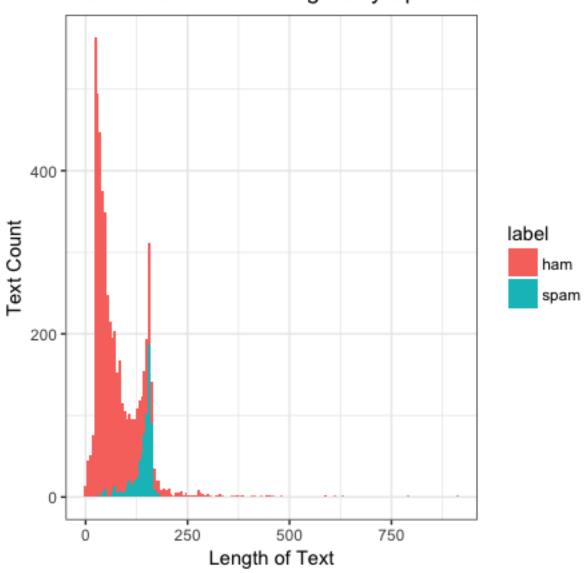# PS6 Li

## Donald Li

## March 2018

# 1 Introduction

The data set I chose to visualize is from Kaggle. Natural Language processing is something that has always interested me. The data set from kaggle contains information over text messages. It then has a label stating if they are "ham" real or spam "fake". My first step was to create column names for the labels. Then I set the label to factors and the texts to characters with the as.factor and as.character.

Going into it I knew I was going to create a bag of words model. However, I also knew the length of texts would be important to determine the classification of the texts. So as to not affect my future model I created a new data set called "messages" to plot the spam and ham messages based off text length.

After that I used the "TM" package to breakdown the texts. I changed all the texts to lower, stemming, removing white space, removing punctuation, removing stop words and remove numbers. At first I tried using the "quanteda" package but I was having issues getting the word cloud visualizations to work.

From there on I created 3 word clouds. One was a word cloud of the entire corpus and the other two are the most common words in each class, ham/spam. The word clouds are important because it can give me clues as to which words are more predictive than others in the classification. Text analytics usually has high dimensionality so being able to reduce noise will help with my predicitions.
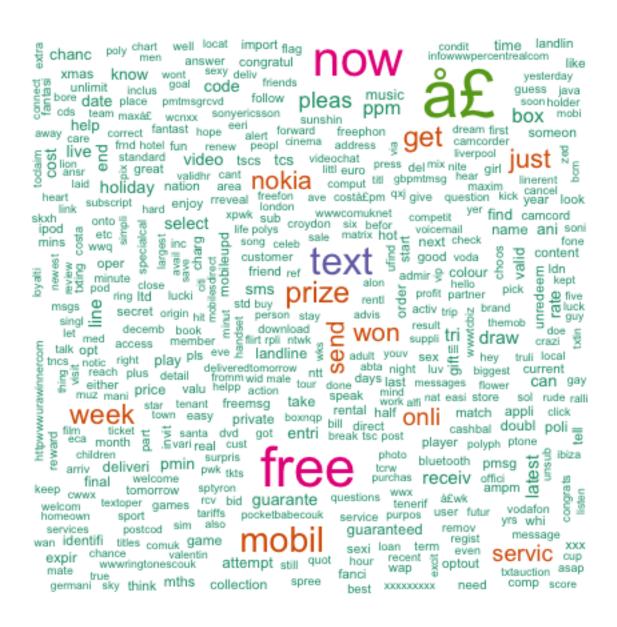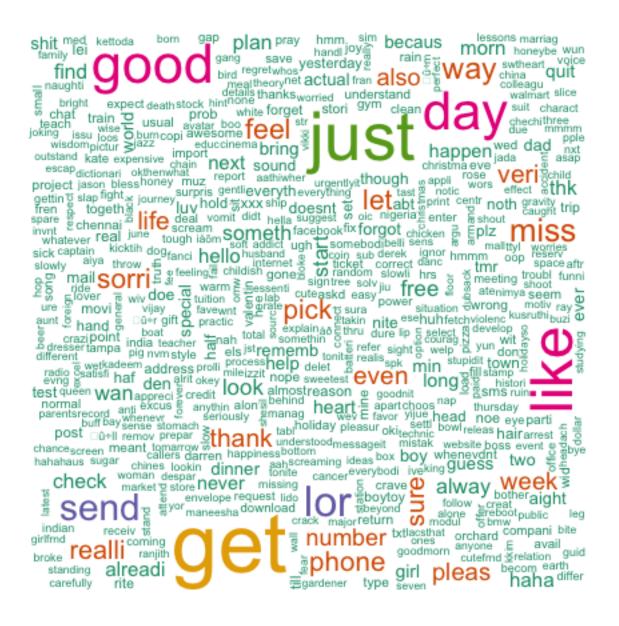
# Distribution of Text Lengths by Spam vs Ham

Figure 2: Ham

Figure 3: Spam