# JHU Data Science Capstone Project - Milestone Report

*Sunday, March 29, 2015*

## 1  Synopsis

This report explores the corpus, HC Corpora, from http://www.corpora.heliohost.org/ in preparation for the John Hopkins Data Science Capstone project. The project is a Shiny app that takes as input a phrase (multiple words) in a text box input and outputs a prediction of the next word. This document explains only the major features identified of the data including summary statistics about the data sets. It also reports any interesting findings that we amassed so far.

In addition, the report elaborates our goals for the eventual app and algorithm and briefly summarize our plans for creating the prediction algorithm and Shiny app. It also serves as a basis for collecting feedback on our plans for creating a prediction algorithm and Shiny app.

## 2  Data Processing

### 2.1  Loading packages

We used different R packages for producing report. More specifically:
For text manipulation, we used: RWeka, stringi
For text mining, we used: SnowballC, tm
For data visualization, we used: ggplot2

### 2.2  Downloading data

The HC Corpora where downloaded from www.corpora.heliohost.org . See the readme file at http://www.corpora.heliohost.org/aboutcorpus.html for details on the corpora available. The files have been language filtered but may still contain some foreign text.

The file downloaded had a large size of 548MB. After unzipping the file, we find the following directories:

```
## [1] "de_DE" "en_US" "fi_FI" "ru_RU"
```

Each of them represent Corpora for a specific language. We will use English Corpora here for our proficiency with the English language. If we list the files in the english corpora directory, we see that it contains files for news, blogs and twitter.

```
## [1] "en_US.blogs.txt"    "en_US.news.txt"    "en_US.twitter.txt"
```

### 2.3  Loading data

When loading all of the english Corpora in the data file, we are actually loading around 151 million words in 2 million lines. That's about four times the size of the Encyclopedia Britannica.
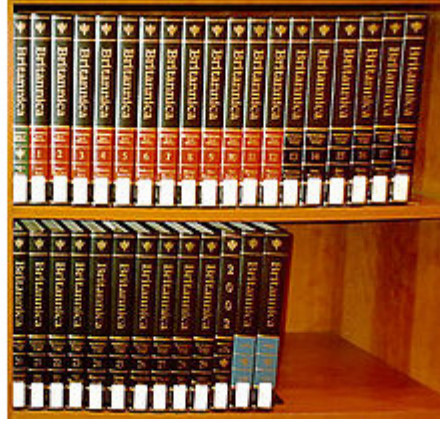
Figure 1: 15th edition of the Britannica. *Wikipedia.*

## 2.4   Data Features and Summary Statistics

To get a better sense of the size of this data, the following summary statistics elaborates the main features of the corpora.

```
##               filename  size_MB total_words    lines mean_words
## 1      en_US.blogs.txt 200.4242    38154238   899288   42.42716
## 2       en_US.news.txt 196.2775    35010782  1010242   34.65584
## 3    en_US.twitter.txt 159.3641     2136398   167155   12.78094
## 4                Total 556.0658    75301418  2076685   34.65582
```
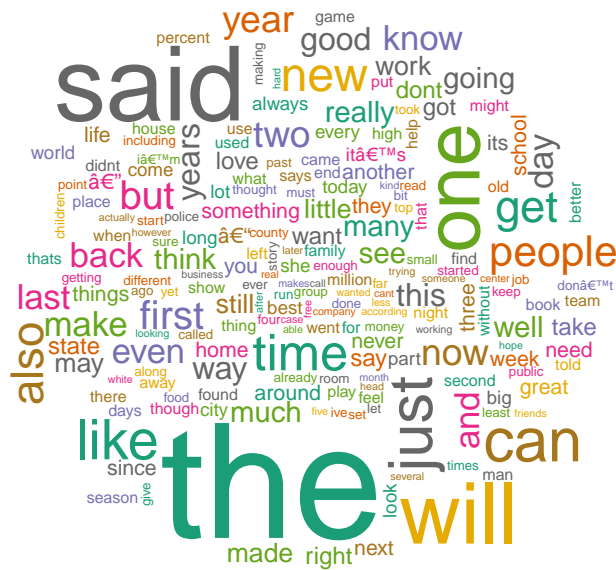
## 2.5   Data Sampling

Beacause the data size is too large which has resulted in many slowdowns and freezes in the computer while exploring the data, I have taken a sample that constitutes 10% of the original data. 10% of each of the files downloaded for the english corpora, the number of lines of the sample is around 200,000 lines verses 2,000,000 lines for the original corpora.

## 2.6   Data Cleaning

Before delving into tokenization of the corpora, i.e. extracting the unique n-gram words or phrases, we will clean the corpora through the following transofrmations:
1- Converting the corpora text into UTF-8 encoding.
2- Removing numbers.
3- Removing unnecessary punctuation.
4- Optionally removing English stopwords, we might cancel this step later as we progress as it may affect the perfromance of our predictive model.
5- Optionally removing profanity words, we might cancel this step later as we progress as it may affect the perfromance of our predictive model. 6- Converting all alphabetical characters into lowercase.
7- Removing extra white space such as tab, etc. . .
8- Stemming the different words, i.e. removing any suffixes such ing, er, s, etc. . .
9- Converting the text into plain text document.

## 2.7 Corpora Tokenization

### 2.7.1 Most Frequent Unigrams (1 word)

#### 2.7.1.1 Bigrams Word Cloud   The word cloud shows the most common unigrams larger in scale.



Figure 2:

#### 2.7.1.2 Unigrams Histogram   The histogram shows the different frequencies of the top unigrams.

### 2.7.2 Most Frequent Bigrams (2 words)

#### 2.7.2.1 Bigrams Word Cloud   The word cloud shows the most common bigrams larger in scale.

#### 2.7.2.2 Bigrams Histogram   The histogram shows the different frequencies of the top bigrams.

### 2.7.3 Most Frequent Trigrams (3 words)

#### 2.7.3.1 Trigrams Word Cloud   The word cloud shows the most common trigrams larger in scale.

#### 2.7.3.2 Trigrams Histogram   The histogram shows the different frequencies of the top trigrams.
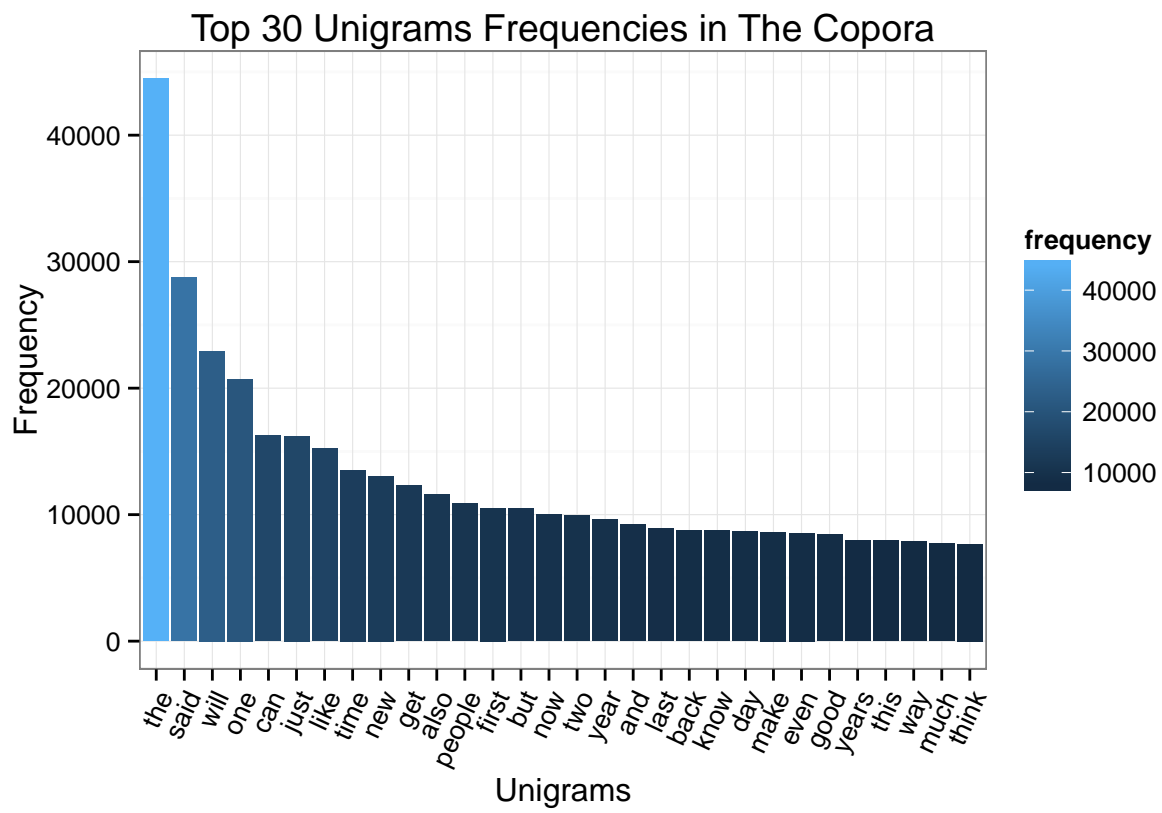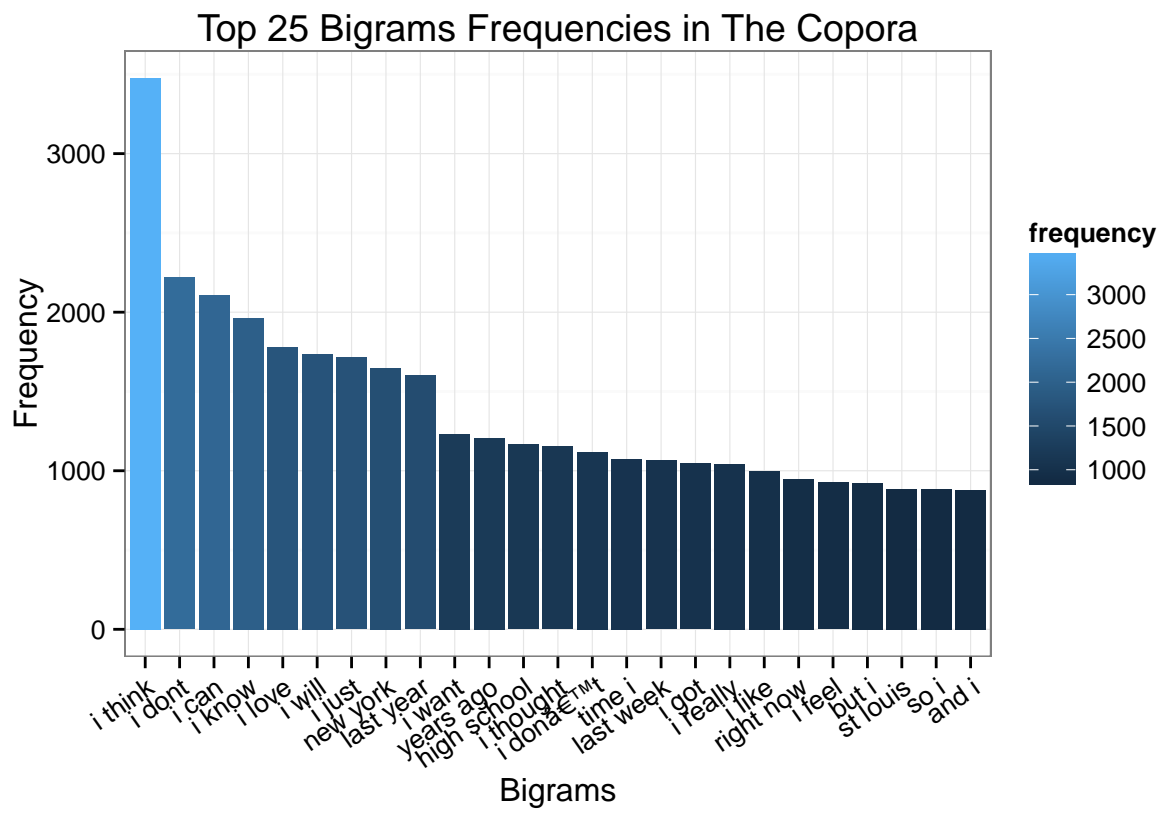
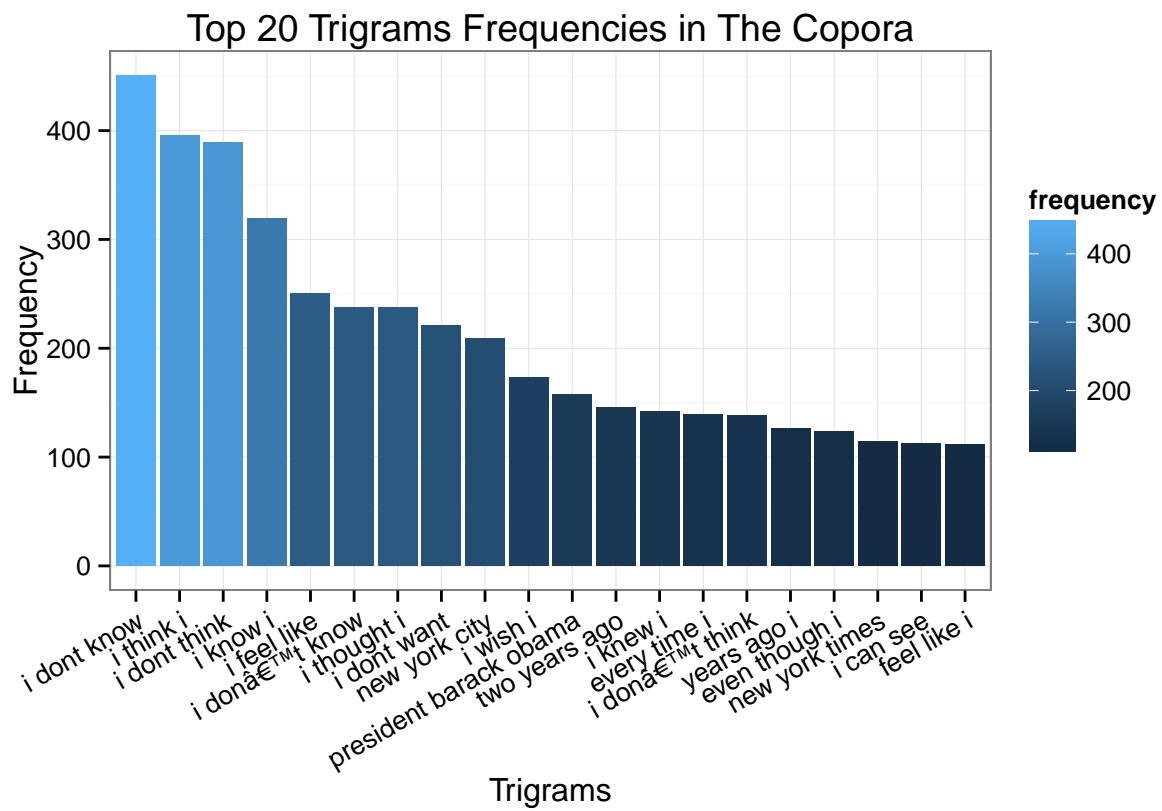Figure 3:

Figure 4:

Figure 5:

Figure 6:

Figure 7:

### 2.7.4 Most Frequent Quadgrams (4 words)

**2.7.4.1 Quadgrams Word Cloud** The word cloud shows the most common quadgrams larger in scale.



Figure 8:

**2.7.4.2 Quadgrams Histogram** The histogram shows the different frequencies of the top quadgrams.

# 3 Shiny App

## 3.1 Text Prediction Strategy

- Preprocessing a specific language corpora using the HC copora or augmenting it with other corpora.

- Building ngram models from 1 and upto maybe 5-7 words depending on the richness of the corpora.

- Building a table for the frequencies of each different ngram model.

- Merging the different tables into one consolidated frequency table for all the ngram models.

- Given a text string of n words, the prediction algorithm will match these words in any of the ngram tokens in the frequencies table, starting with the n+1 gram tokens, if it does not match any n+1 gram token, the last n-1 words of the input string is matched in the n gram tokens this time and the process
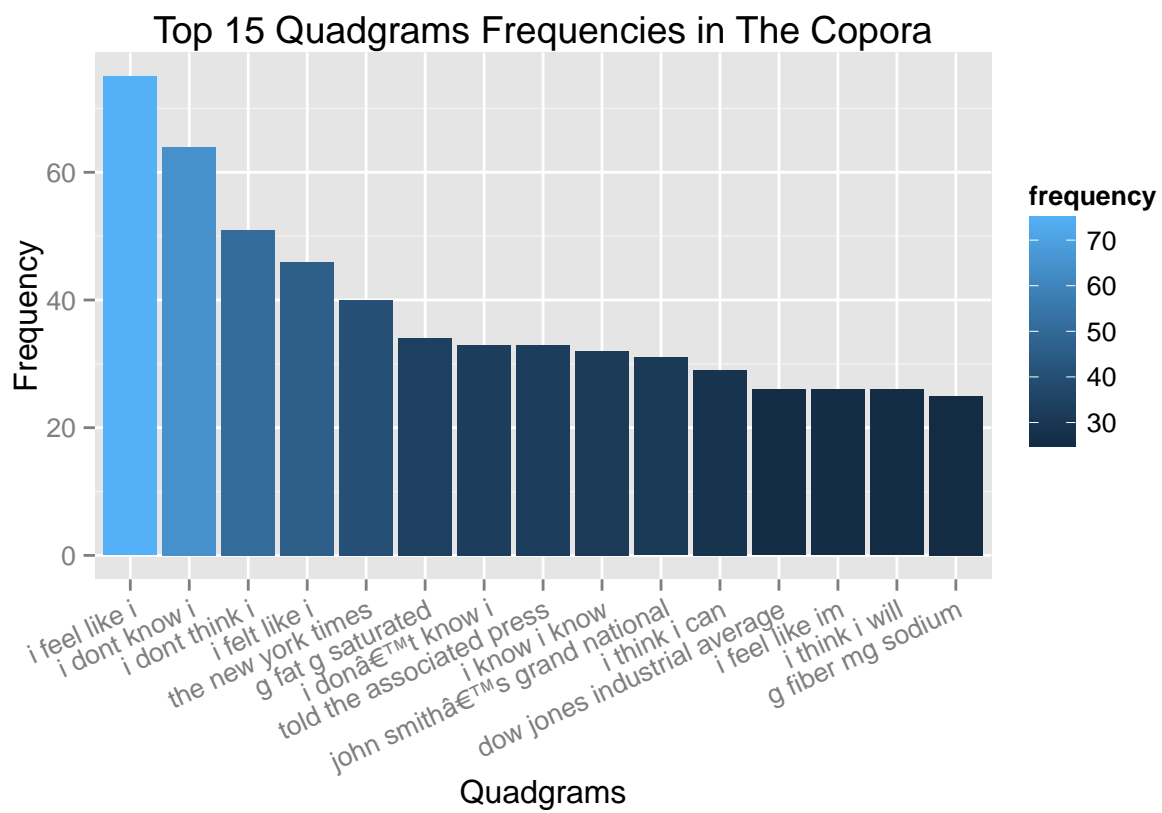
Figure 9:

continues until a match is found or the last word, i.e. 1 word of the input text is matched in one of unigram tokens.

- If input is matched, the last word of the n+1 gram tokens with the highest frequencies are recommended for the user to complete his/her input text.

## 3.2 Shiny App Plan

- Deciding on the optional cleaning steps to include in the final prediction model. For example, will the final model include stopwords and profanity words?

- Looking for augmentary english corpora that can be used to enhance the prediction to be built.

- Completing the n-gram model, I have already built the unigram, bigram, trigram, and quadgram models. However, I might be considering ngram models based on 5-7 words as well if the dataset provided is rich enough to enable me to extract these models.

- Deisgning and prototyping the user interface for the shiny app.

- Implementing the Shiny App.

- Documenting the App.
- Developing a short pitch for the Shiny app.

# 4 Conclusion

This document reports the main data features of the HC Copora downloaded from http://www.corpora.heliohost.org/. For each text file in the English copora, the report shows the size in addition to other summary statistics such as, wordcount, number of lines, and the mean number of words per line. After building the unigram, bigram, trigram, and quadgram models, the report shows the most frequent n-gram words/tokens along with their frequency in the English Corpora. Finally, the report elaborates the text prediction strategy and the plan for the development of the Shiny App, the final product.