

Milestone Report Submission

Rodrigo Bertollo de Alexandre

Thursday, September 18, 2014

As requirement for the first report submission, we are supposed to demonstrate what we have done so far. At first I've download the [Capstone Dataset](#) under the **Announcements** page within the *Task 0 - Understanding the problem*. After extraction, I started to work with the *en_US* folder containing the three files: - *en_US.blogs.txt* - *en_US.news.txt* - *en_US.twitter.txt*. Since these files are a large amount of text data, the best approach method is the **tm package**. It is important to highlight that inside the *DirSource* argument under the *Corpus* function, the encoding must be set to “**UTF-8**” to avoid encoding errors.

The file representing *blogs* data has about **nine hundred thousand lines**, the *news* has about **seventy seven thousand lines** and *twitter* files has almost **two million and four hundred thousand lines**.

Before starting to work with these files, some formatting was done. At first, words and letters were transformed to lower case. Due to the presence of strange characters the best way to transform all words to lowercase was with the use of the **stringi** package. Moreover, later on this package also allowed to do further modifications to the text and *tokenization*.

Since we are working with the English language, all lowercase i and i' where transformed back to I and I' respectively. Unfortunately, the code that I had done for multiple substitutions was slower than doing them separately.

The function *removePunctuations* was not used due to the fact that this function would also remove . and ' from the text. Without the dots different phrases would be wrongly linked together and without the apostrophe some contractions words would lose its sense, example “*I’m*” to “*Im*”; “*it’s*” to “*its*”. Therefore, all end phrasing punctuations where transformed to dots, and all other punctuations with exception of apostrophe where deleted.

Further on, other modifications like removing extra white spaces, numbers and separating every dot by a new line was also executed. The separation of new lines for every dot was important because during the use of n-grams different phrases wouldn't be linked as one.

Because bad languages use (*cursing and swearwords*) also make a critical an important part of the language, they were not removed of the database, but they will be removed during the prediction model.

For later use and backup purposes the file was constantly saved with the function `save(us_files, file="us_files.RData")`. The data was removed from the Corpus with the function `as.list` and converted to a list document file. A matrix was constituted using the argument `unlist` and `byrow = T`. All matrix lines that had less than three words were deleted from the matrix, and saved with the function `save(clean_data, file="clean_data.RData")`. The word count was performed with the **qdap package**.

At the end, the cleaning and processing resulted in a matrix containing 7644814 lines. Tokenization was performed with the **stringi** package.

The tokenization resulted in an absurd number of words (509478 words). This happened because some entries that were not words were also considered as words.

##		Var1	Freq
## 100	aaaaaaaaaiwwll		1
## 99		aaaaaaaaages	1
## 98		aaaaaaaaayyyyyyyyyyyyyyyyy	1
## 97		aaaaaaaaahhhhhhhhhhhhhhhhhh	1
## 96		aaaaaaaaaargh	1
## 95		aaaaaaaaaaaaammmmmmmmm	1

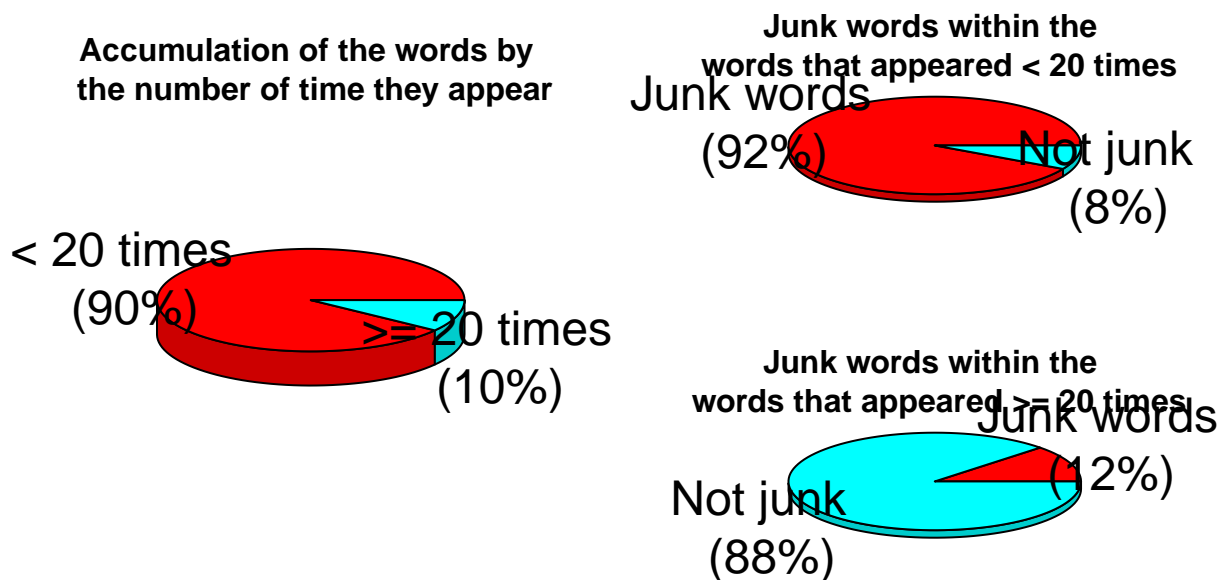
Because of this, I've created English word list containing about 655000 english words gathered from innumerable databases containing nouns, verbs, names, names of cities and countries, swearwords and much more.

I used this file to filter existing words from my *data_dic*. In this way I was able to create a "not word dictionary" for later use to exclude the lines containing these junk words from my ngrams.

Interestingly, because of this procedure I was able to identify some writing mistakes that should be considered of correction before using the prediction models.

```
##      Var1  Freq
## 258615   lol 46614
## 220335  it's 38814
## 210141   im 29989
## 121276 don't 29421
```

Furthermore, if among the 509478 token words, about 458278 appear at most 19 times, whereas 51200 had a frequency equal or higher than 20 times. After the removal of the existing words (as mentioned before), about 83% of the overall were consider junk, (92% of the less than 20 times and 12% of the ≥ 20 times).



.pdf

This junk words will only be deleted after the creation of the n-grams because only the parts that are affected by it will be removed and not the whole phrase.

The n-grams will be performed by dividing the *clean_data* file in 8 different parts due to memory RAM usage. For this, the **ngram package** will be used. tryCatch function is used for avoiding errors due to phrases with less than the minimal words during the n-gram loop. The same code will be used for all 8 files and for all different n-grams.

After this, there will be further processing and data cleaning, and separation of the n-grams by alphabetical order in different files.