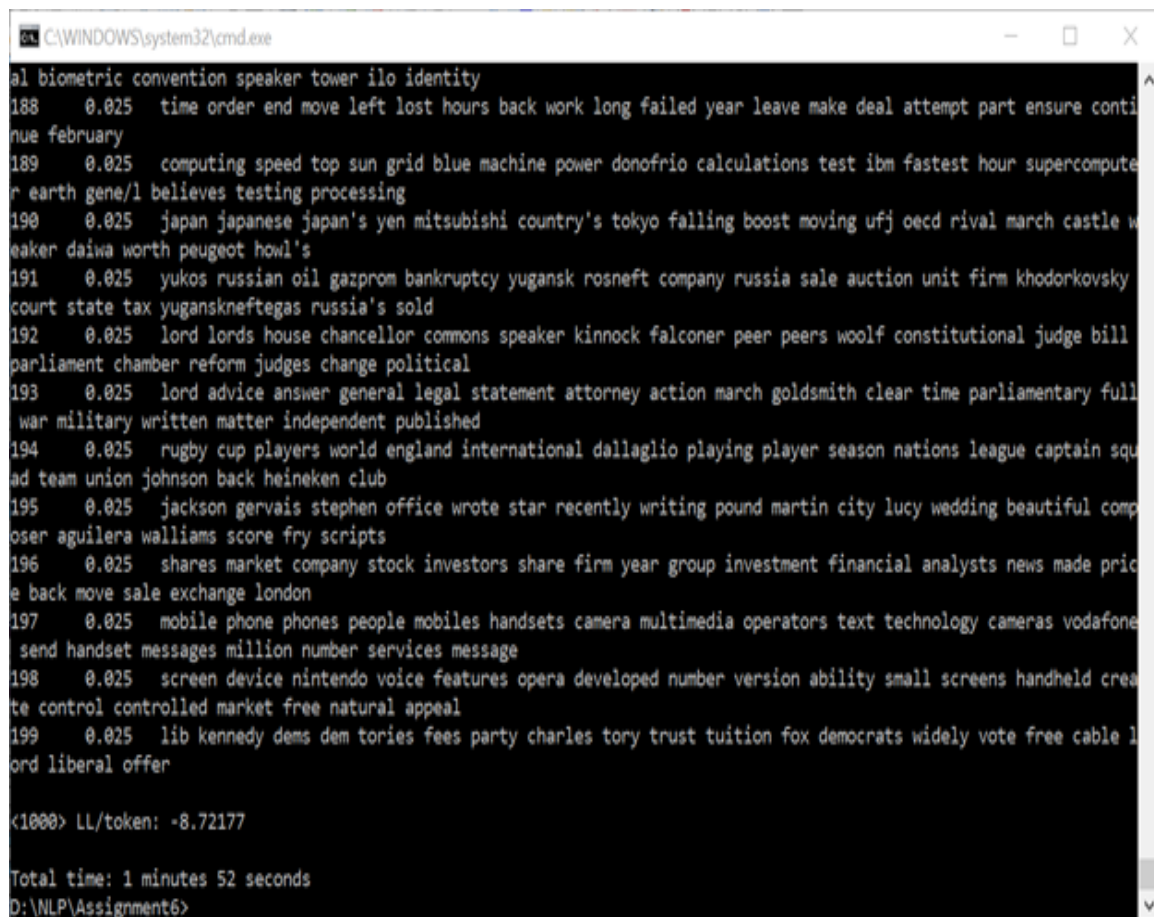ADVANCED NATURAL LANGUAGE PROCESSING

Homework 6

Anusha Ramamurthy

a. Pick some topic modeling software, e.g., one of the following, and get it working:

We initially chose R to do this task. Unfortunately we had trouble installing the library tm and slam. Mallet seemed to be more user friendly.

b. Pretend you are a social scientist or a humanities scholar and come up with some rough question you want to investigate. Gather a small amount of appropriate data and run the topic modeling software on that data, to learn the models.

We decided to use BBC articles for the year 2004-05. The question we had in mind was to identify the top 200 topics published by BBC for this year. The data was split into different genres, but we decided to run mallet on the entire set. The below is a screenshot of the keywords identified for each topic.

c. Investigate the output of the tool and give a brief assessment. Do the topics (i.e., the most frequent words) make sense? Can you see the documents annotated in a way that helps answer your question? Do any mistakes in the output connect with what you know about how the system is working? If you were a non-computer-savvy social scientist, what would have been the most significant barriers to using the software (if any)?

We chose the command line argument option to save the outputs in file format. The bbc_200.txt gives us the keywords associated with different the top 200 topics.We found 194 occurances of "football" , 402 occurances of America,  334 occurances of India and 17 occurances of "racial" in these 200 topics.

It definitely is a good way to identify the top trends of the year. It's also in my opinion, a good way to identify the bias if there are any. As "church" appears 48 times among 200 topics. If we ran Mallet over a national daily like an Indian Newspaper we would probably have more news about  topics around say "temples"

The most glaring mistake is that these words may appear many times in the same article/topic, but one would assume that its the most spoken about word. The mallet picks words at random and the results are based on the order. If we jumble up the words we will not get the same results.  If an article uses a lot of matephors and similis to explain the topic, in such cases Mallet would not be able to identify this. For example if he uses dog as an example to illustrate loyalty, Mallet might assume that dog is one of the topics. So it definitely depends on the writing style used.

The command line would definitely be intimidating to a non computer savvy person. Also the outputs are not very clear. Understanding and evaluating the results/outputs produced needs more study.

Acknowledgements
Worked in collaboration with Anup Bharadwaj.

References:
http://mlg.ucd.ie/datasets/bbc.html
https://cran.r-project.org/web/packages/topicmodels/index.html
 (Links to an external site.)

Citation:
D. Greene and P. Cunningham. "Practical Solutions to the Problem of Diagonal Dominance in Kernel Document Clustering", Proc. ICML 2006.

Copyright:
All rights, including copyright, in the content of the original articles are owned by the BBC.
Contact Derek Greene <derek.greene@ucd.ie> for further information.
http://mlg.ucd.ie/datasets/bbc.html