

Home Assignment 3

Javier Galindos Vicente, 214373IV

The student would like to present the work offline in class.

I. TEXT MINING

In this section, some speeches of President Biden have been selected as the text dataset. The first step is to preprocess the dataset. In particular, the punctuation symbols, numbers, capitalization and the words ('ll, 're, 've, the, and) have been removed towards having a more precise analysis. The histogram of the most common words is presented in Fig. 1 and the word-cloud with the most frequent words is presented in Fig. 2. One could note observing this analysis that the dataset is composed of speeches from an American politician due to present of words like america, people, nation, country, democracy, united, etc.

The documents are represented in a Document-Term Matrix (DTM) for simplicity. Before clustering the data, it is necessary to remove sparse terms. Specifically, the maximal allowed sparsity is 0.15 in this dataset. The similarity measure used to compute the distances is the cosine similarity for every algorithm due to the better performance than other metrics (i.e. euclidean, manhattan).

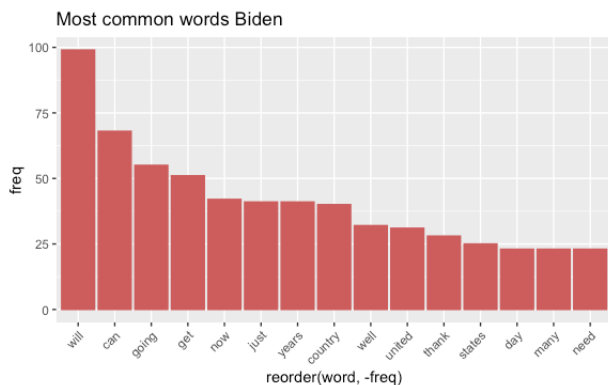


Fig. 1: Histogram of the text dataset.

To cluster the dataset, several techniques such as k-means, hierarchical clustering or k-medoids have been performed. Other probabilistic algorithms like EM or density based methods such as HDBSCAN have been explored, but the results are not optimal. The three clusters generated by the aforementioned algorithms agree in having 3 centroids between the most common words. One of them is formed from key words where Pres. Biden wants to persuade some ideas with words like state or troops. Another centroid depicts words that he uses to engage his followers with examples (e.g. god, bless, thank, work, today, etc.) The last centroid represents a bigger cluster with the frequent words Biden mentions to differentiate himself and

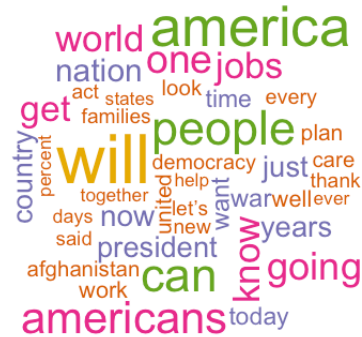


Fig. 2: Word-cloud of the text dataset.

set a political line with words (e.g. states, american, united, together, nation, people, protect, families, country, etc.).

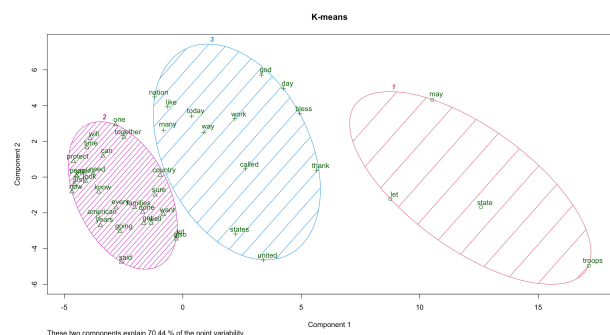


Fig. 3: K-means clustering.

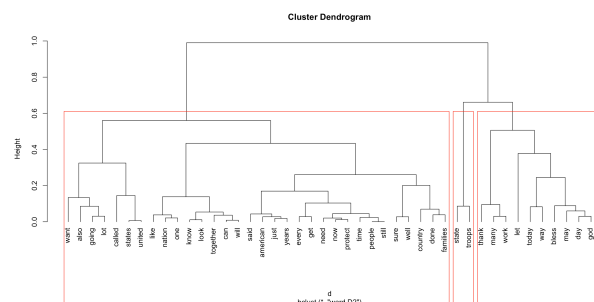


Fig. 4: Hierarchical clustering

Numerically, both the centroids and the decision boundaries could be explained in term of two components function of the frequencies (see Fig. 3 and Fig. 5). These components explain 70.33 % of the point variability. In particular, for the k-means clustering, it could be observed that the component 1 is almost capable of depict a decision boundary, whether the component

