

MOTIVATION

Online discussions about public policy are essential to good policy making. By examining public discourse, opinions and perspectives can be acknowledged and considered to find adequate solutions. Such a public policy involves energy policies, more specifically, the Energy East pipeline project announced by TransCanada on August 1, 2013. This project was emblematic, since it became heavily debated on social platforms like Twitter, among multiple groups on economic, political, environmental, and moral grounds. This project aims to categorize the users to identify different demographic's positions on this topic.

DATA PROCESSING

The data processing pipeline that was followed and the set of tools used for this project as seen in Figure 2 were the following:

1. Data Collection

Twarc2 is an application used to collect Twitter posts from users' timelines and users' profile description, retrieved from the official Twitter API in CSV format for easier data manipulation.

2. Data Wrangling

Pandas was used to select, filter, and clean users' respective tweets through dataframes to convert such data into words' vector through the GloVe model. By using a words' vector, composed of a direction and a value, words can be compared to one another. See Figure 1.

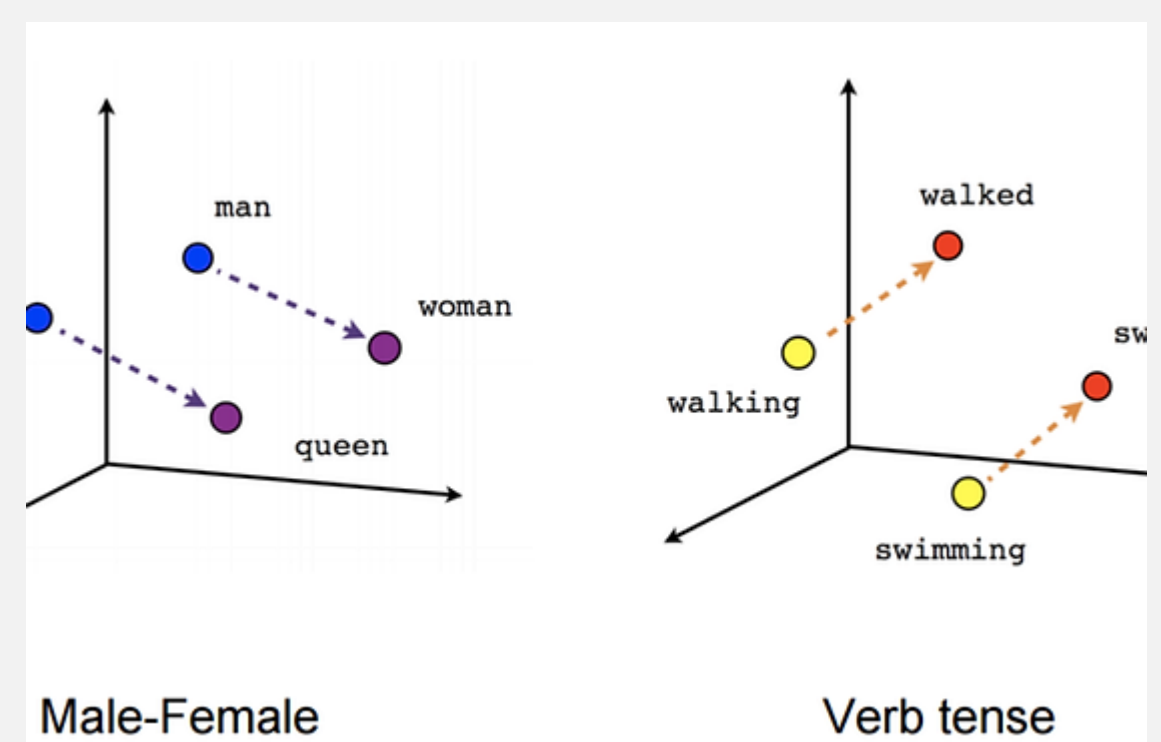


Figure 1: Visual Example of Word Vectors. Taken from: <https://www.synthesisproject.org/resources>

3. Data Analysis

An algorithm called agglomerative clustering was selected to compare the sum of vectors for each tweet and pair groups of similar data together, to later explore such clusters' similarities to assign them a label as a category.

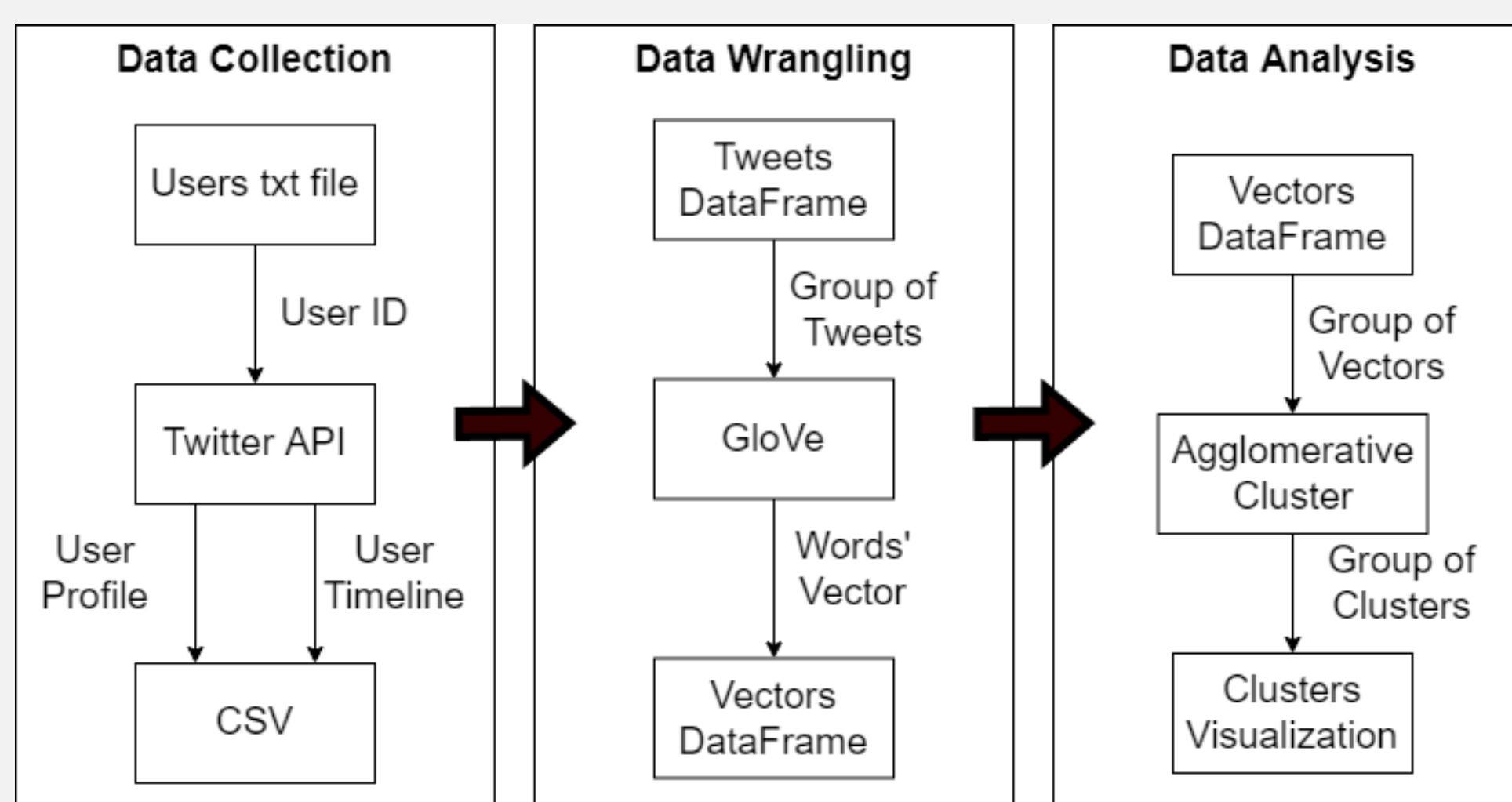


Figure 2: The Data Processing Pipeline.

METHODS & ANALYSIS

Words' Vector and GloVe

A 50-dimensional embedding of words was used to obtain the word's corresponding vector.

A word has a vector, which is a direction and value. Since GloVe has pre-trained vectors, each word has a pre-existing position, which on a 2-dimensional graph, these words would have an (x, y) coordinate. But since they are 50-dimensional embeddings, each word has fifty numbers describing their position instead of two. See figure 3. By summing the vector of each word in a tweet, and then summing the tweet's vector for 50 tweets of each user, which will be called the total vector, the overall significance of what the user engages with is captured. See figure 4 and 5 above for total vectors.

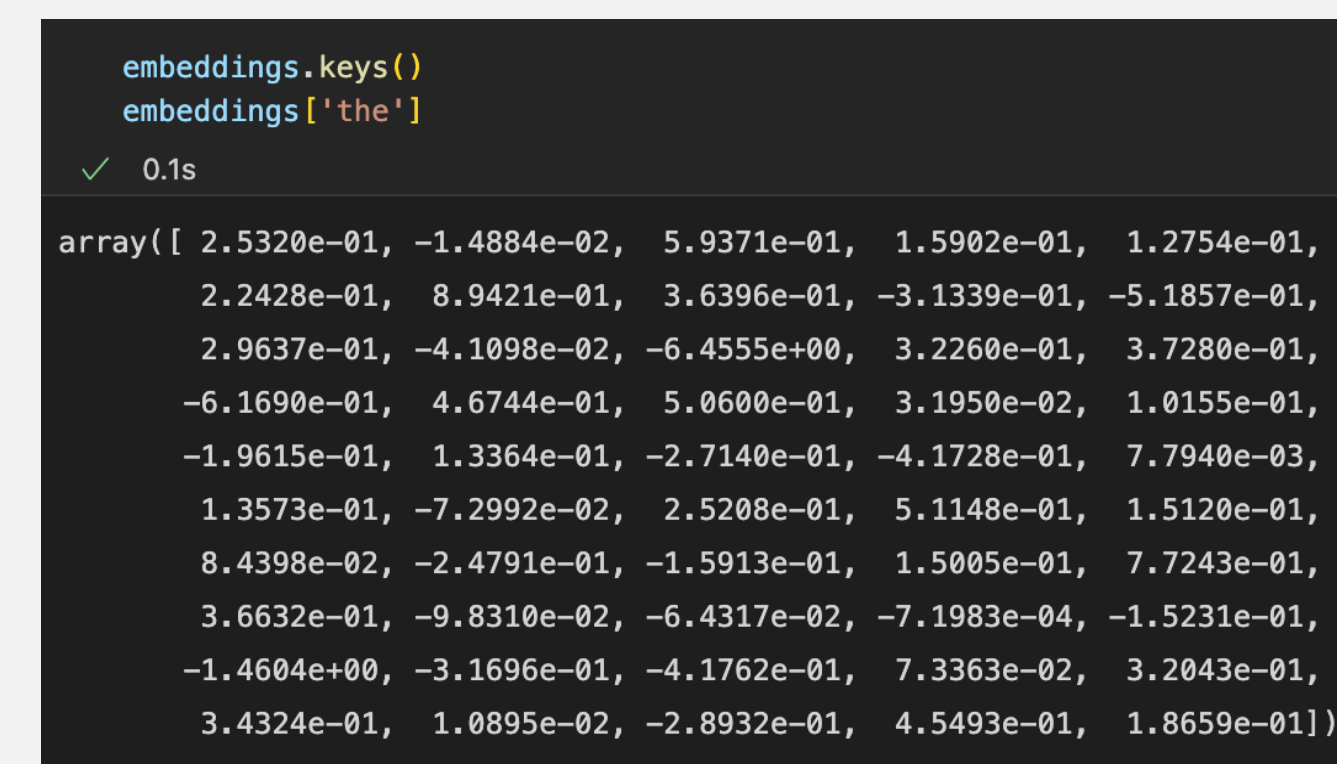


Figure 3: vector for the word "the"

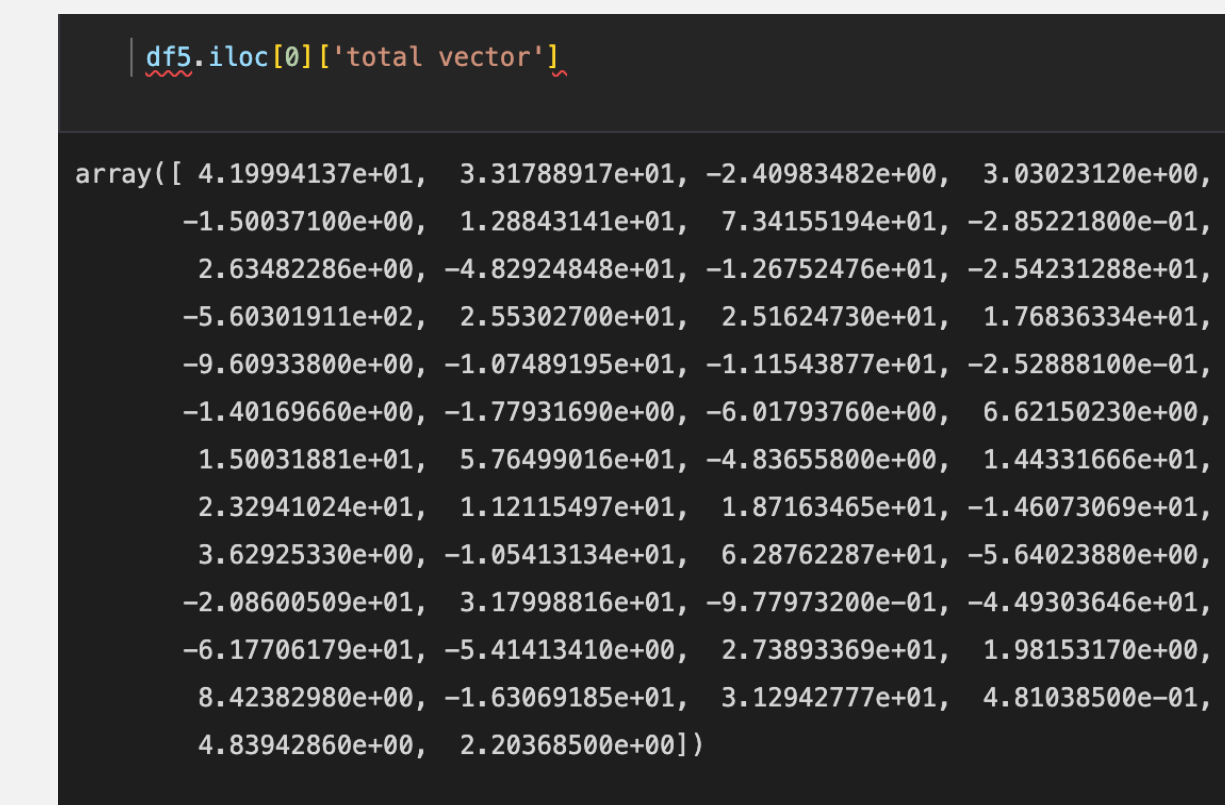


Figure 4: Total vector for user 0

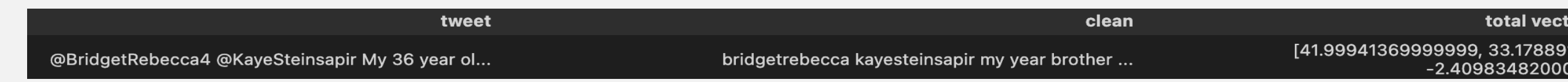


Figure 5: user 0's tweets, cleaned tweets, and total vector

Agglomerative Clustering

Agglomerative clustering is an algorithm used to group data. By allowing the computer to group similar total vectors, the groups can be later studied and analyzed for similarities, and then categorized. For example, below is a graph of three different clusters. Each cluster represents a group of tweets that have similar overall vectors. Referring to figure 6, the red points could be twitter accounts with total vectors closer to politics, while the blue points could be twitter accounts with total vectors closer to environment, and the green points could be total vectors closer to economics.

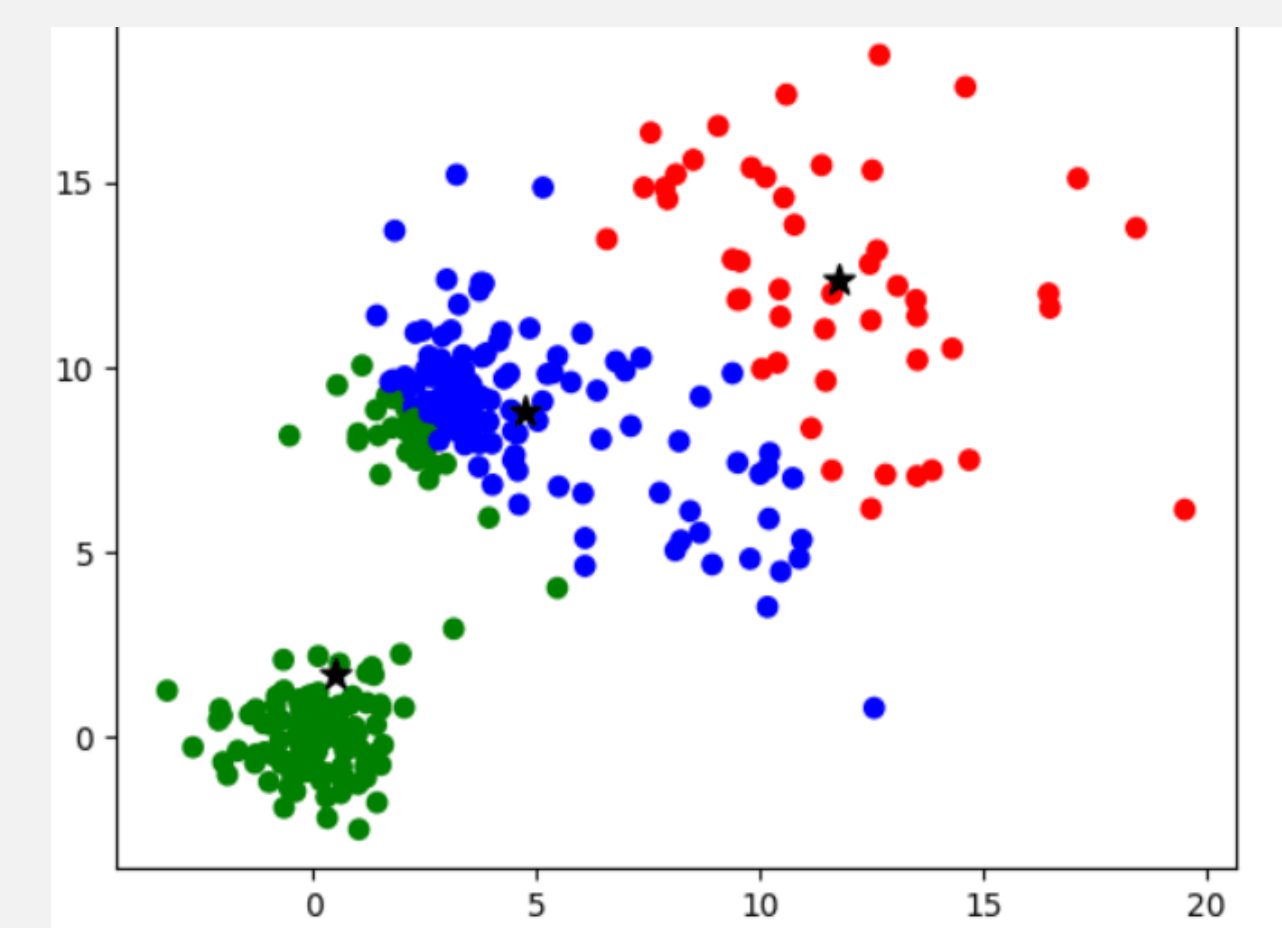


Figure 6: an example of agglomerative clustering
Taken from: <https://dev.to/piyushbagani15/k-means-clustering-25hf>



Figure 7: 5 Twitter users from each cluster

Analysis

The 30 000 Twitter users' latest 50 tweets were extracted, and then cleaned with the program NLTK so only the nouns were left. Only the nouns were analyzed since words like "I", "had", or "is" do not have any significance to the users' overall engagement. Then, the 50 tweets of each user were processed in GloVe for a total vector. After inputting all the users' total vectors in the clustering algorithm, and inputting 8 as the desired number of groups, the computer created a graph with 8 clusters, see figure 8. Each cluster is a group of similar total vectors. To find the optimal number of clusters for a data set, the elbow method is used, see figure 10.

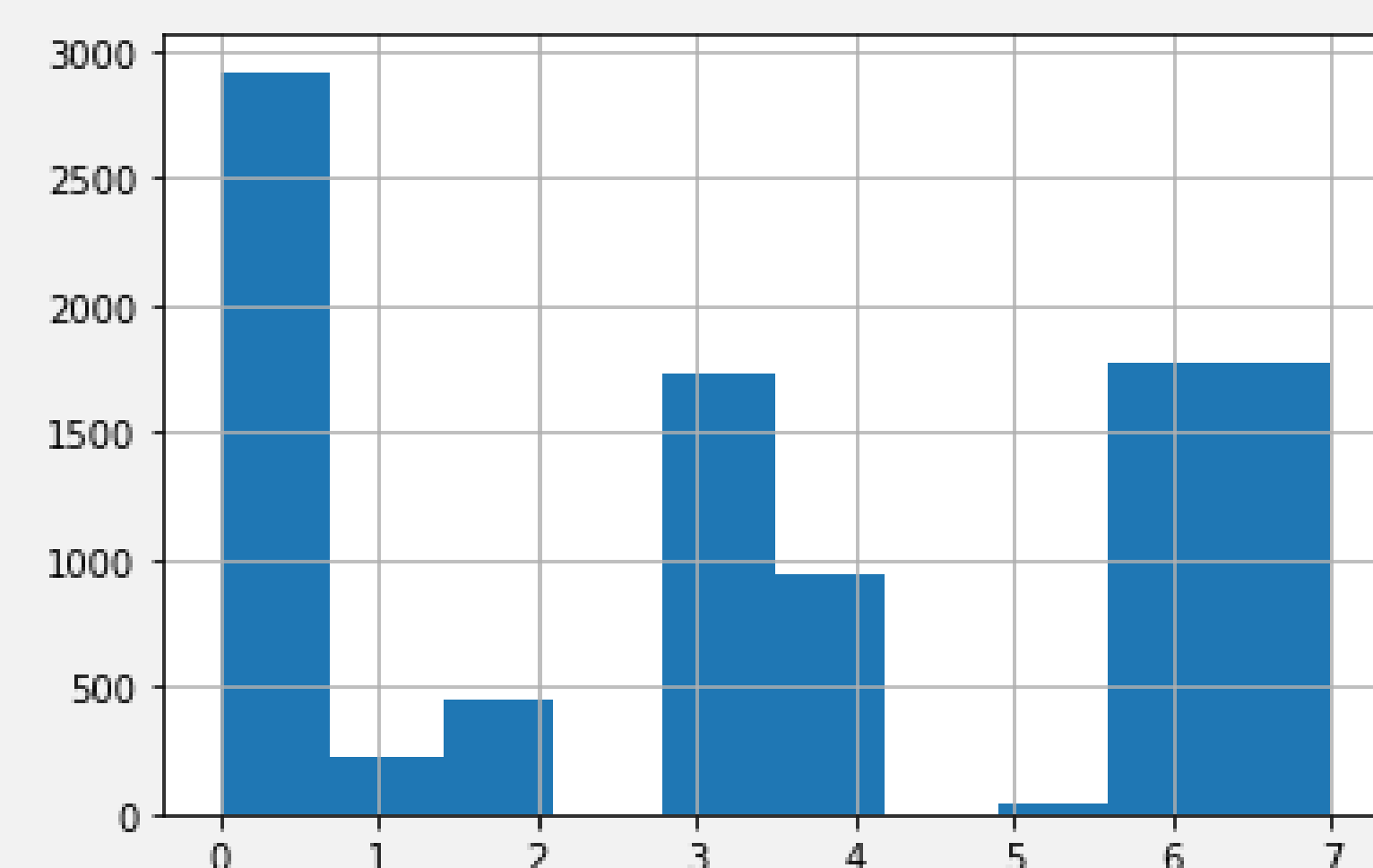


Figure 8: Graph of the 8 clusters

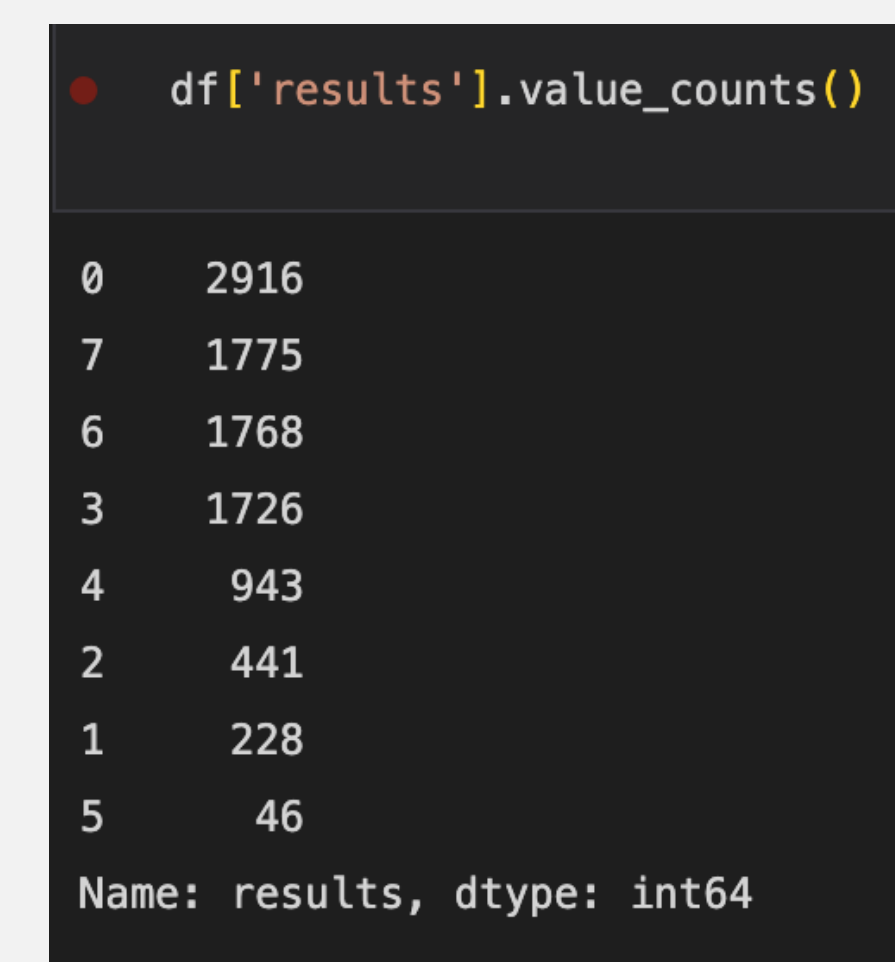


Figure 9: The number of Twitter users in each cluster

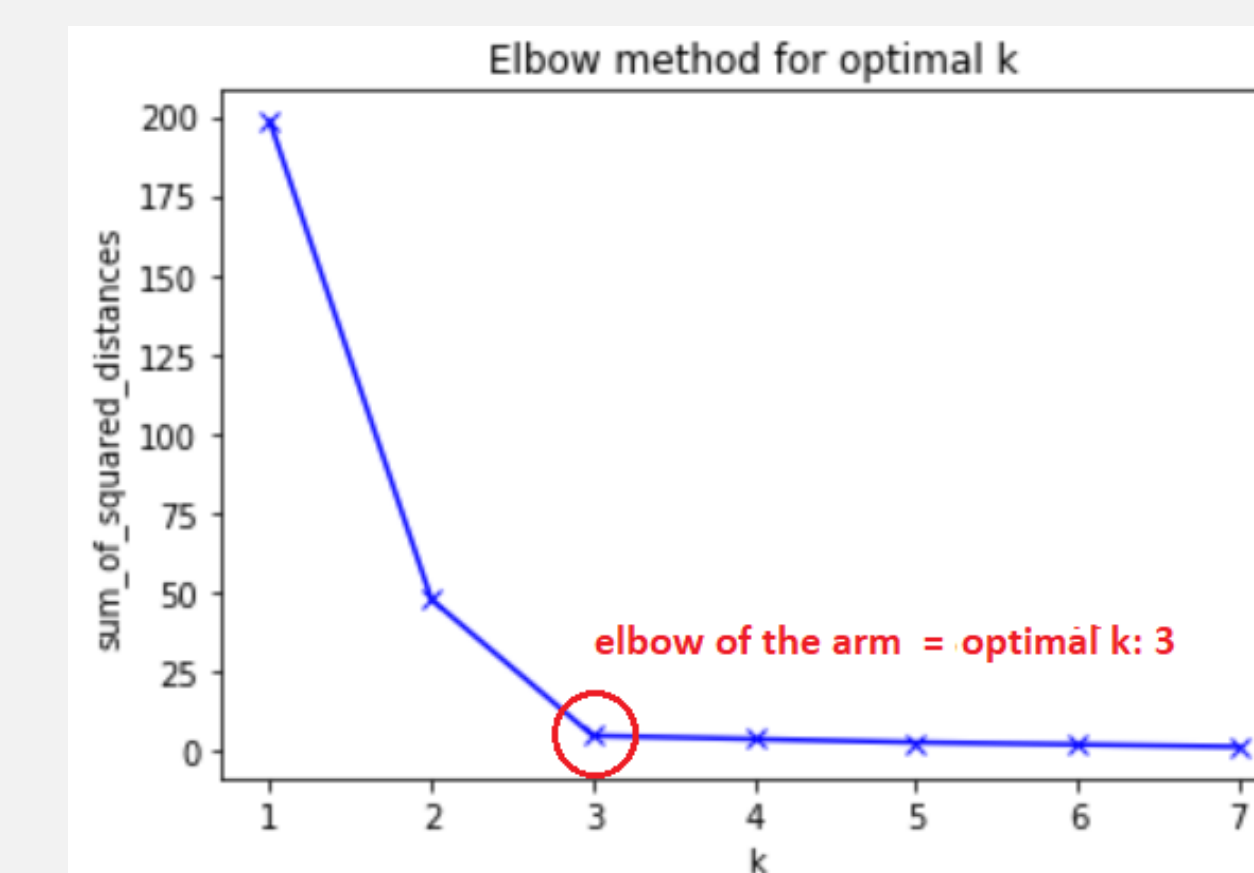


Figure 10: the sum of squared distances as a function of the number of clusters
Taken from: <https://stackoverflow.com/questions/82443970/finding-the-optimal-number-of-clusters-using-the-elbow-method-and-k-means-clust>

RESULTS

Face Validity

After examining some user accounts from each cluster, some similarities were found. In some cases, these similarities were obvious, but for the most part, there was a lot of overlap. Below are some generalizations made after examining 5 users from each cluster.

- Cluster 0: anti-Trudeau, right wing
- Cluster 1: French text
- Cluster 2: climate change activists
- Cluster 3: left wing
- Cluster 4: economy activists
- Cluster 5: Spanish text
- Cluster 6: miscellaneous
- Cluster 7: not clear

The only two clusters that were strongly consistent were cluster 0, cluster 1, and cluster 5. The other clusters were moderately consistent, while cluster 7 did not have a clear theme.

When analyzing the tweets in each cluster, tweets were grouped together based on the object of conversation, but not the sentiment. For example, a user who tweeted: "I hate Trudeau", and another who tweeted: "Trudeau is great", were placed in the same cluster since the noun of both sentences is "Trudeau." This made it hard to find consistent similarities in each cluster.

Additionally, some of the tweets were image-based, to which the user only tweeted emojis or basic sentiments as text. This was not properly analyzed by GloVe since, as stated above, GloVe cannot capture sentiment. This too made the data hard to successfully cluster.

CONCLUSIONS

By using the GloVe program and the clustering algorithm, the 30 000 Twitter users were categorized into 8 groups. By having a model that can categorize users based on the meaning of words, the perspectives of each group can be studied for better policy making. This fits into the bigger project as it allows for demographic trends in persuasion and manipulation tactics to be studied.

Next Steps

The 8 clusters were moderately consistent, however, some clusters had a lot of overlap. To mitigate this, sentiment analysis should be used to cluster users not only on the topic discussed, but also on their position towards it. Additionally, more investigation is required to identify the ideological affinities present in each cluster. Lastly, more calculations on the optimal number of clusters is required.

ACKNOWLEDGEMENTS

Thank you Dr. Lefsrud, Dr. Stroulia, Mr. Roberts, and Mr. Gutierrez for helping me with this project, and thank you AI4Society and Dr. Lefsrud for sponsoring this project! I would also like to thank everyone else who helped me along the way, and finally, I would like to thank the WISEST team.