

New York Times Articles about “Immigration” Analysis

Zhecan Wang

Overview

Based on the New York Times API, I make a script to stream all the articles that mention the keyword “Immigration” from 1850 to 2005. As a result, I collect around 30000 articles totally. However, in order to better compare this data set with other bills and budget data in time series, I would focus on the articles in the range of 1900 to 2005 (around 15000 totally).

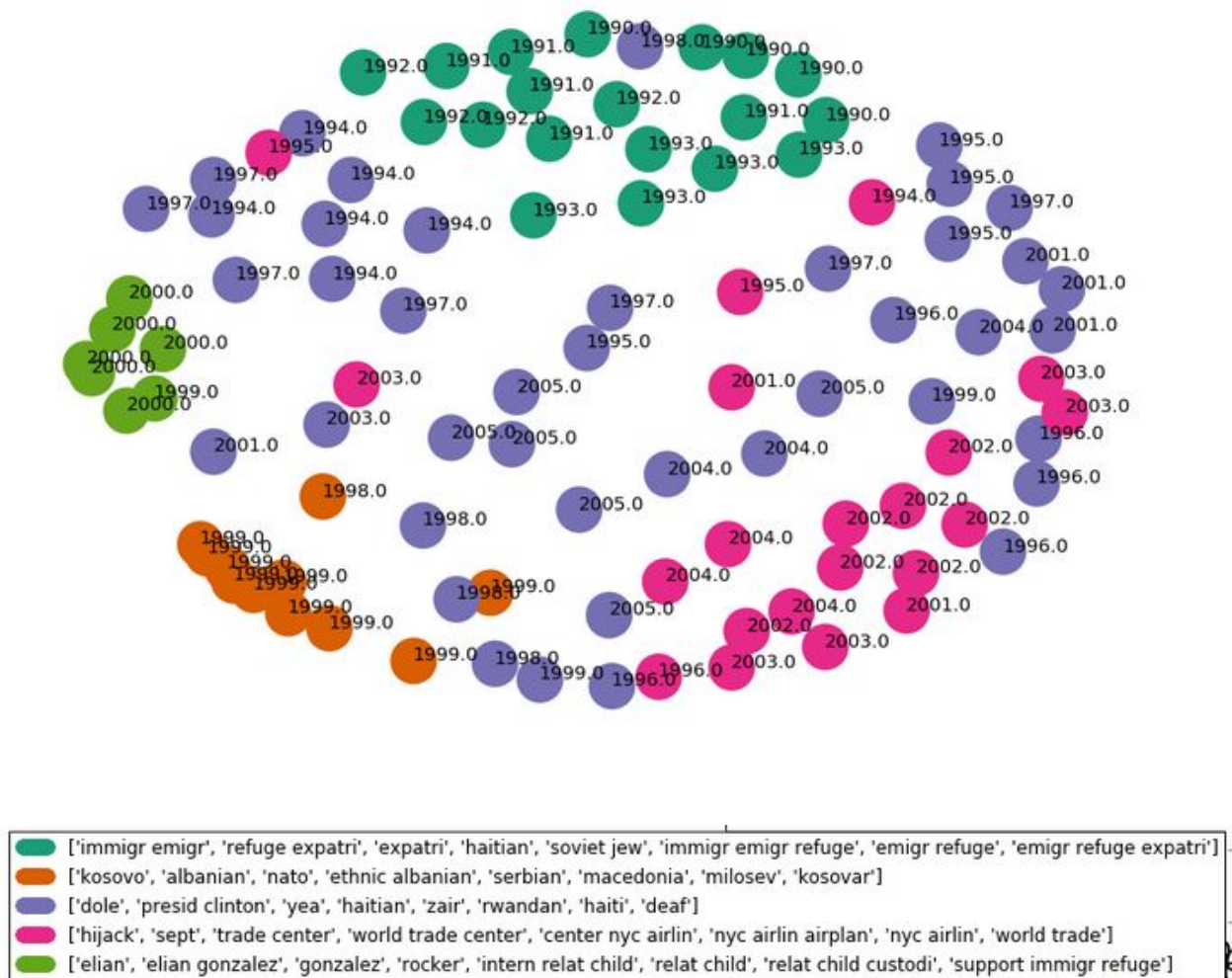
Clustering

I use TF-IDF to analyze all the articles and extract a feature matrix from the data set. In this matrix, each row represents either a group of articles or an article and each column is a feature word. After feeding this matrix in the K-means clustering algorithm, we end up with 5 clusters.

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

In the following picture, I group the articles together into 100 bins among the 15000 articles, so each circle represents around 150 articles. Each circle is originally a feature vector and should be graphed in a high dimensional space but here I project (smash) them all into a two dimensional plane by using the **Multidimensional Scaling** (reduction).

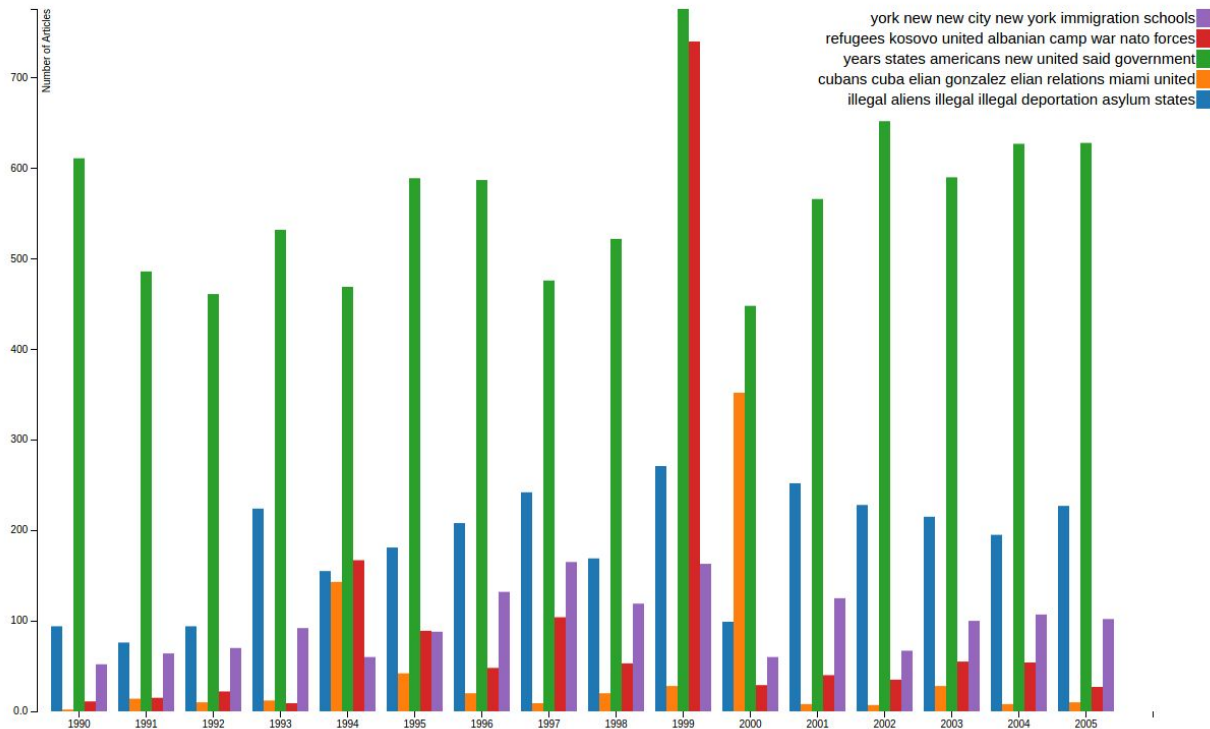
As you may observe, the circles in a cluster tend to be closer together in spatial distance. Also, the year of circles in a cluster also tend to be closer in chronological distance. The clusters' keywords are very distinct from each other. Overall, it is a reasonable clustering.



The following is the overall clusters distribution across years.

Among the five clusters, there are two which may be more interesting. The first one is the “kosovo” (second in order) cluster that tends to center around 1999. The background context can be referred to the Kosovo War, Kosovo Albanians and so on. The second one is the “cubans” (fourth) cluster that tends to center around 2000. The background context can be referred to the Elian Gonzalez custody battle, former Cuban

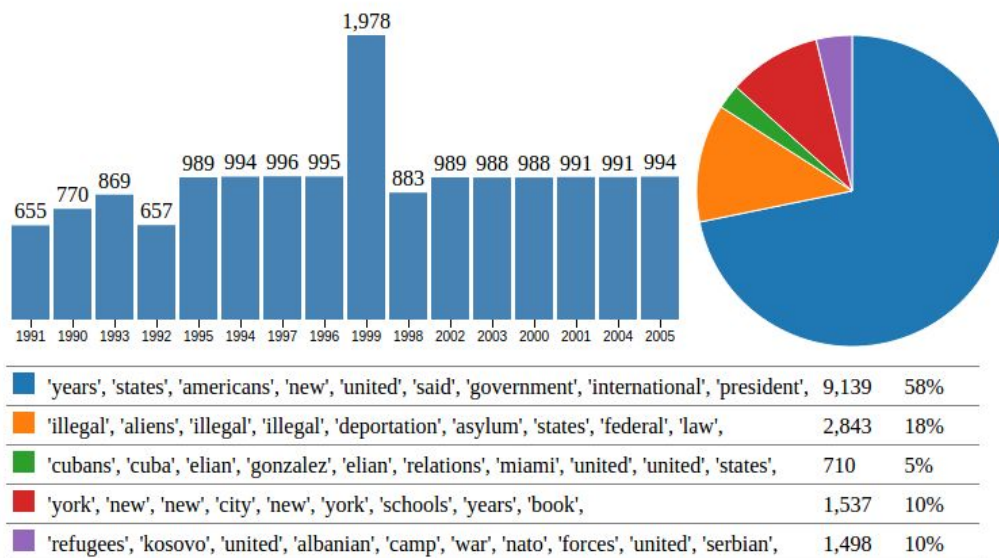
castaway Elián González and so on. The other clusters tend to be more general comparatively.



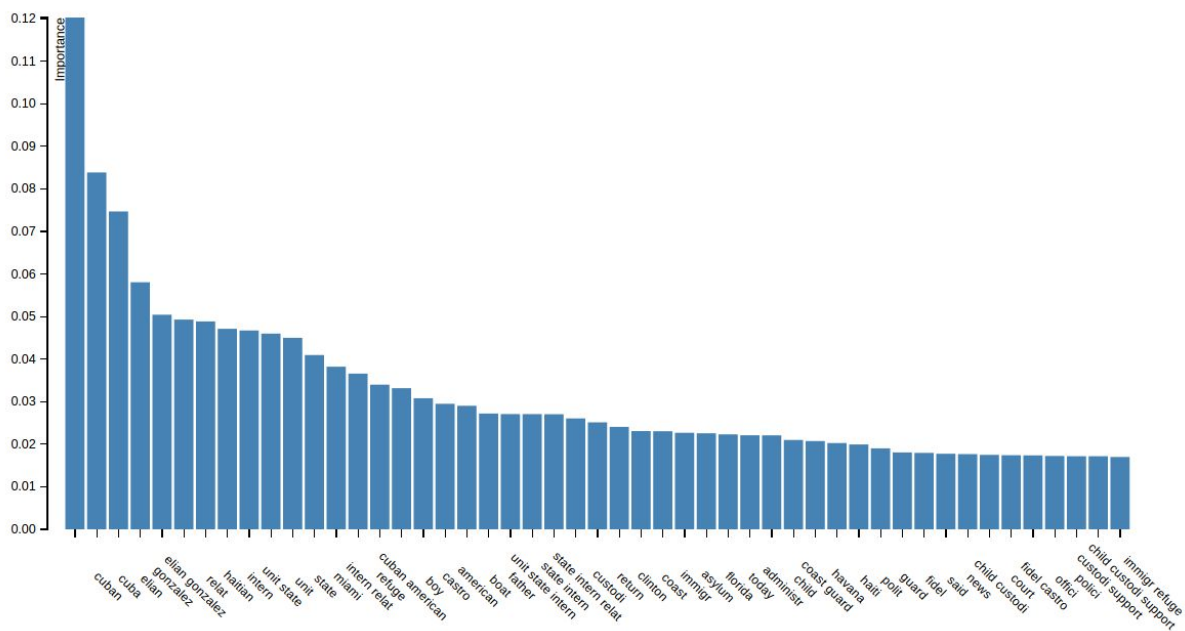
After clustering, we grab the center words from each cluster along with their importance score (from TF-IDF) and use word cloud to visualize. The following picture shows the important words in the “kosovo” (second in order) cluster. The interesting words here are “kosovo”, “albanian”, “serbian” and so on.

Visualization

The following is another D3 visualization I created for the clustering year distribution (The actual interactive script can be found in the github repository).



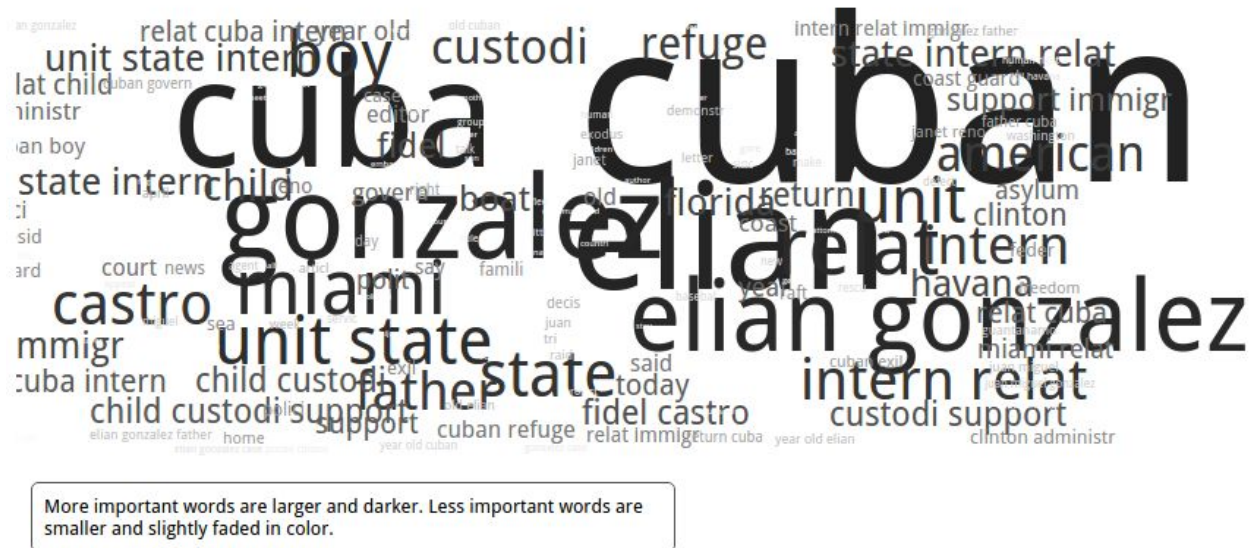
The following is a D3 histogram of the important words in the second cluster by TF-IDF.



More important words are larger and darker. Less important words are smaller and slightly faded in color.

Term	Importance
cuban	0.120
cuba	0.083
elien	0.074
elien gonzalez	0.058
gonzalez	0.050
relat	0.049
italian	0.048
italian	0.046
unit state	0.045
unit state	0.041
state	0.038
miami	0.037
intern relat	0.034
cuban american	0.033
boy	0.031
castro	0.029
american	0.029
boat	0.027
unit state	0.027
father	0.027
unit state	0.027
state intern relat	0.026
state intern relat	0.025
custodi	0.023
return	0.023
clinton	0.023
coast	0.023
immigr	0.022
asylum	0.022
florida	0.022
today	0.022
adminisr	0.022
coast guard	0.021
child	0.021
havana	0.020
haiti	0.020
polit	0.019
guard	0.018
fidel	0.018
said	0.018
child news	0.018
court	0.018
fidel castro	0.018
offici	0.018
polici	0.018
custodi support	0.018
child custodi support	0.018
immigr relat	0.018

As you can see from the word cloud, the words, “cuban”, “elian gonzalez”, “miami” and “custodi support” show up a lot.



Overall, the clustering tells us two important events, the Kosovo War around 1999 and the Elian Gonzalez custody battle around 2000.