

# 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

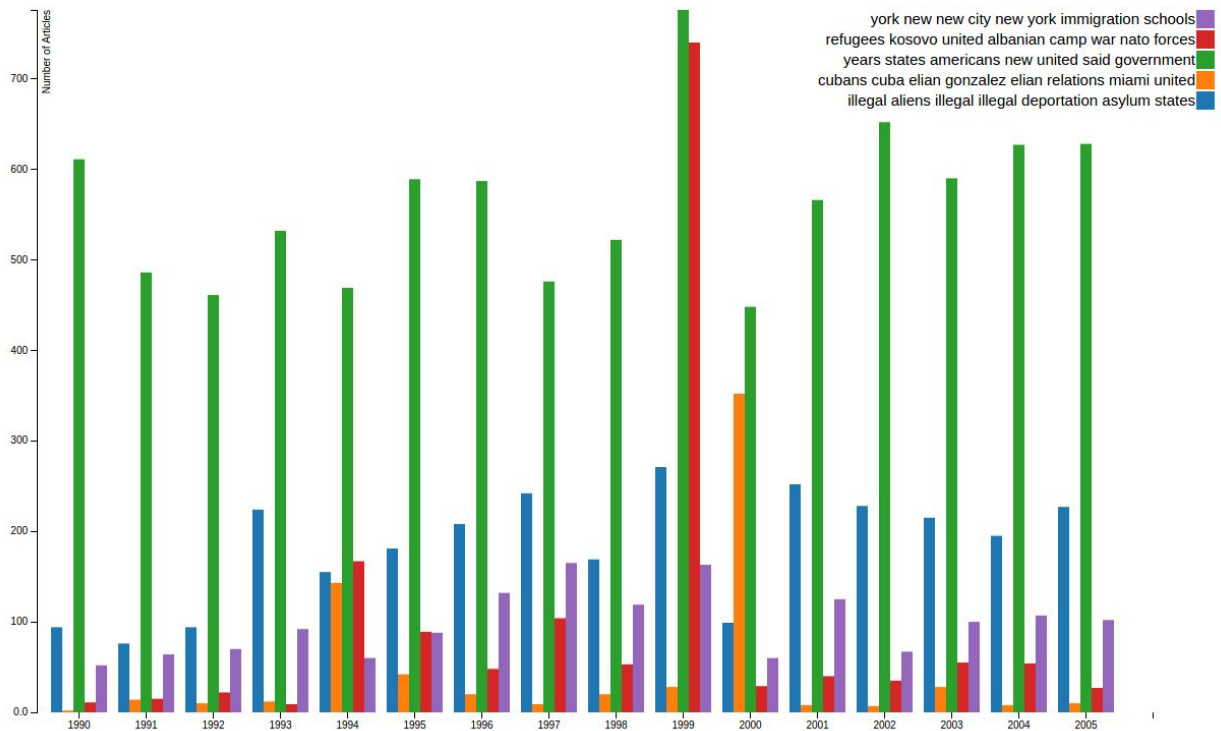
In the following picture, I group the articles together into 100 bins among the 15000 articles, so each circle represents around 150 articles. Each circle is originally a feature vector and should be graphed in a high dimensional space but here I project (smash) them all into a two dimensional plane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



er

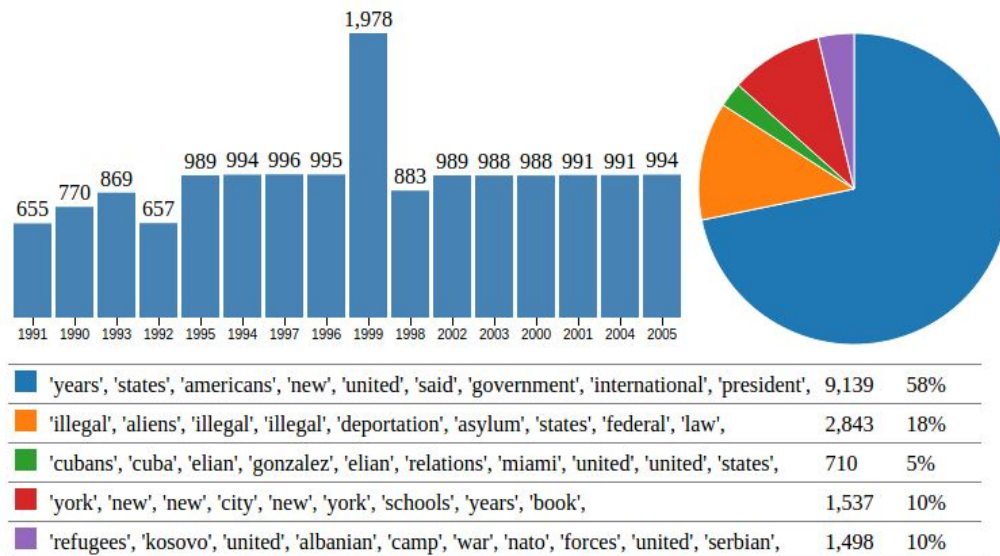
m

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.



## 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).

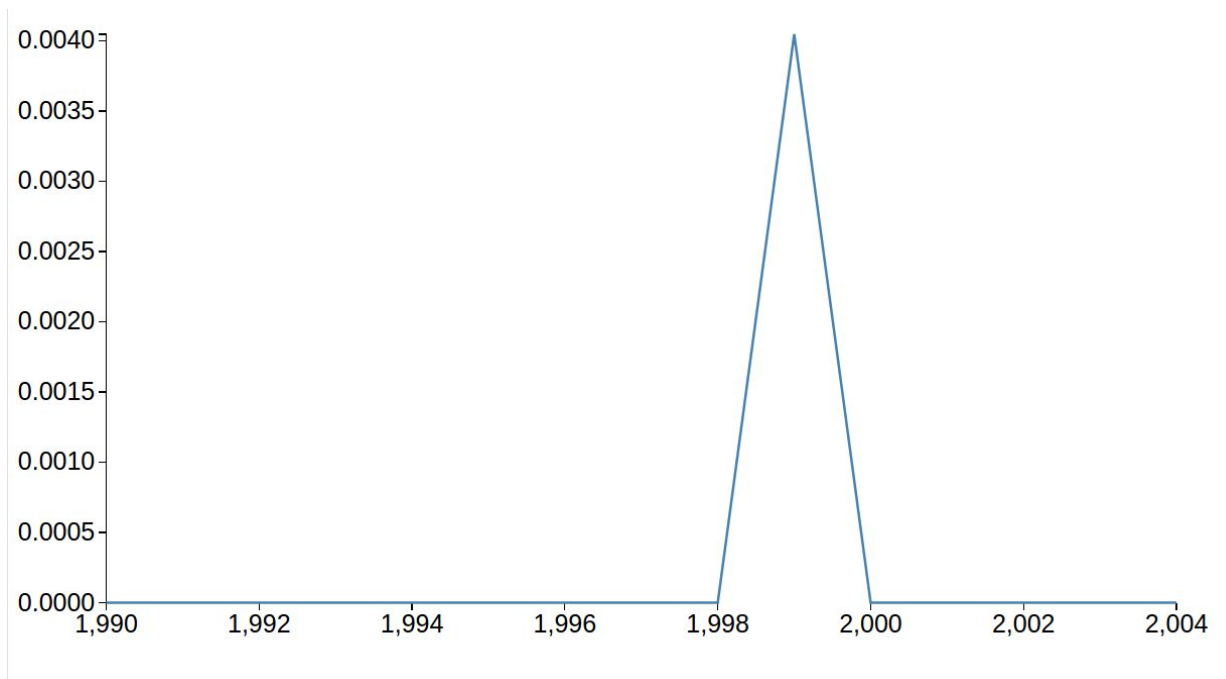


After clustering, we grab the center words from each cluster along with their importance score (from TF-IDF) and use histograms and word clouds to visualize. The following picture shows the important words in the “kosovo” (second in order) cluster. The interesting words here are “kosovo”, “albanian”,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

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

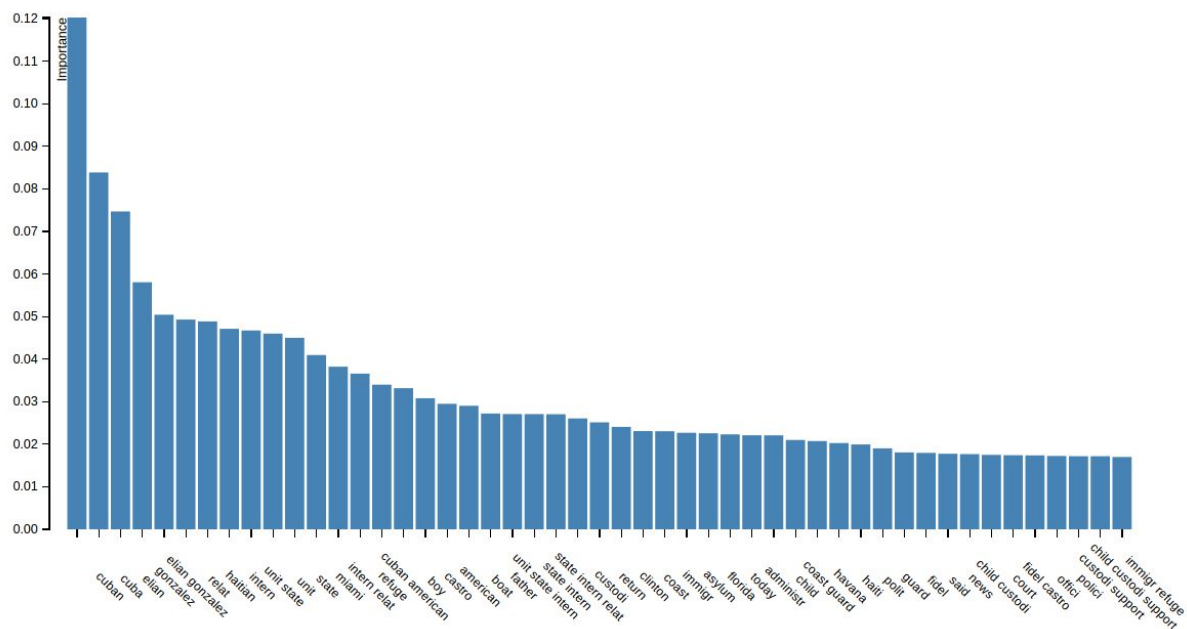


This graph shows the possibility that you can see the word "kosovo" in the new york times "immigration" related articles. The X axis is the years and the Y axis is the normalized possibility

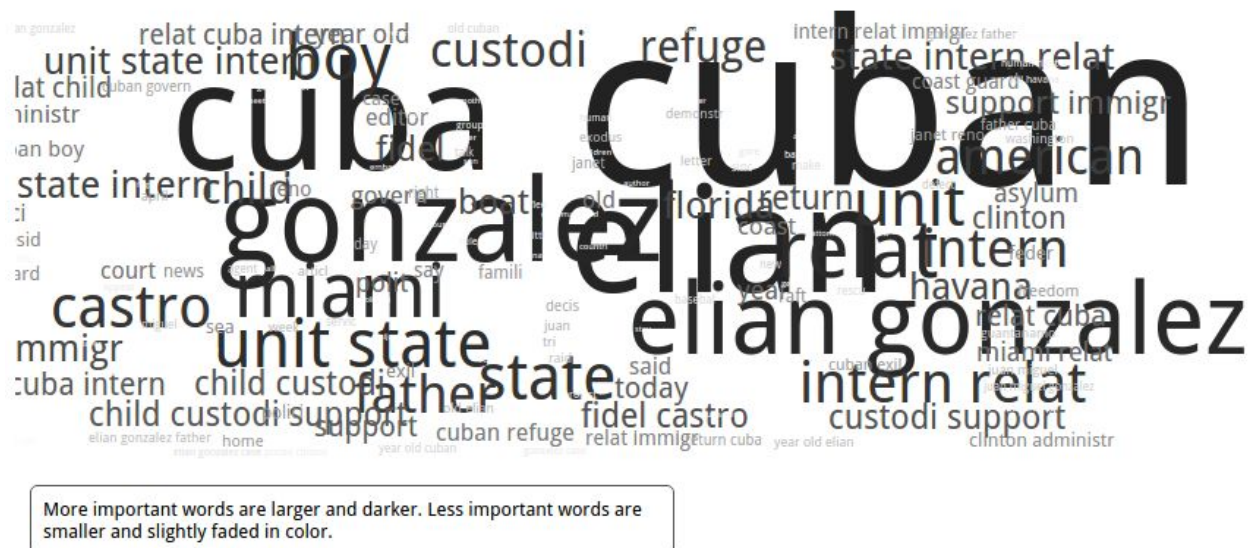




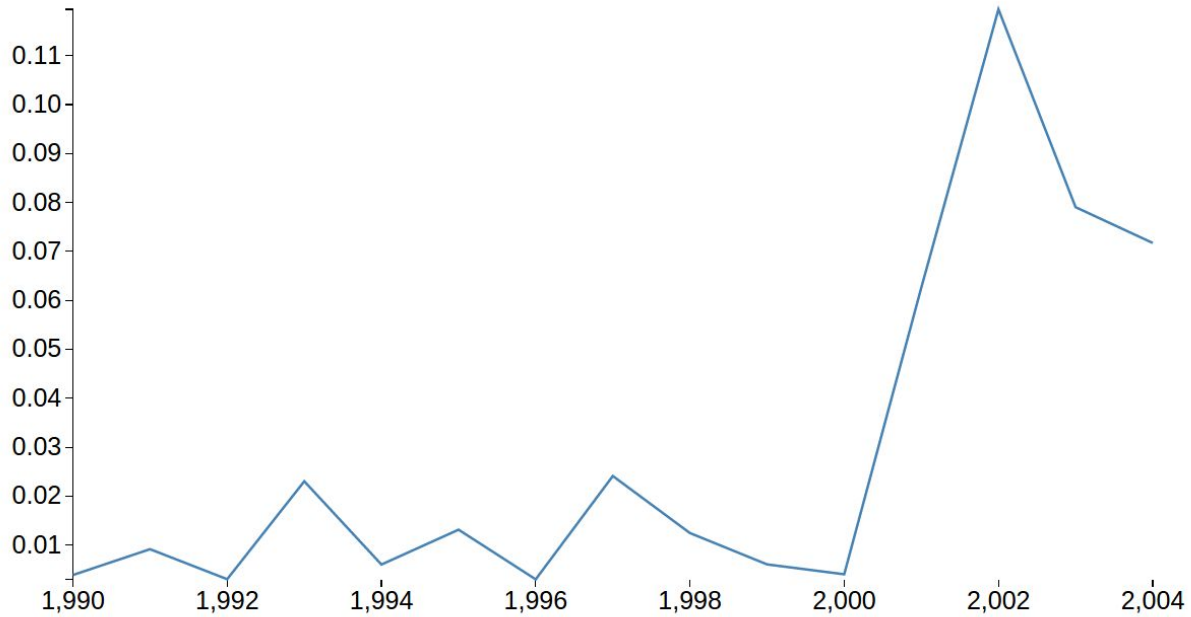
The following is a D3 histogram of the important words in the last cluster by TF-IDF



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

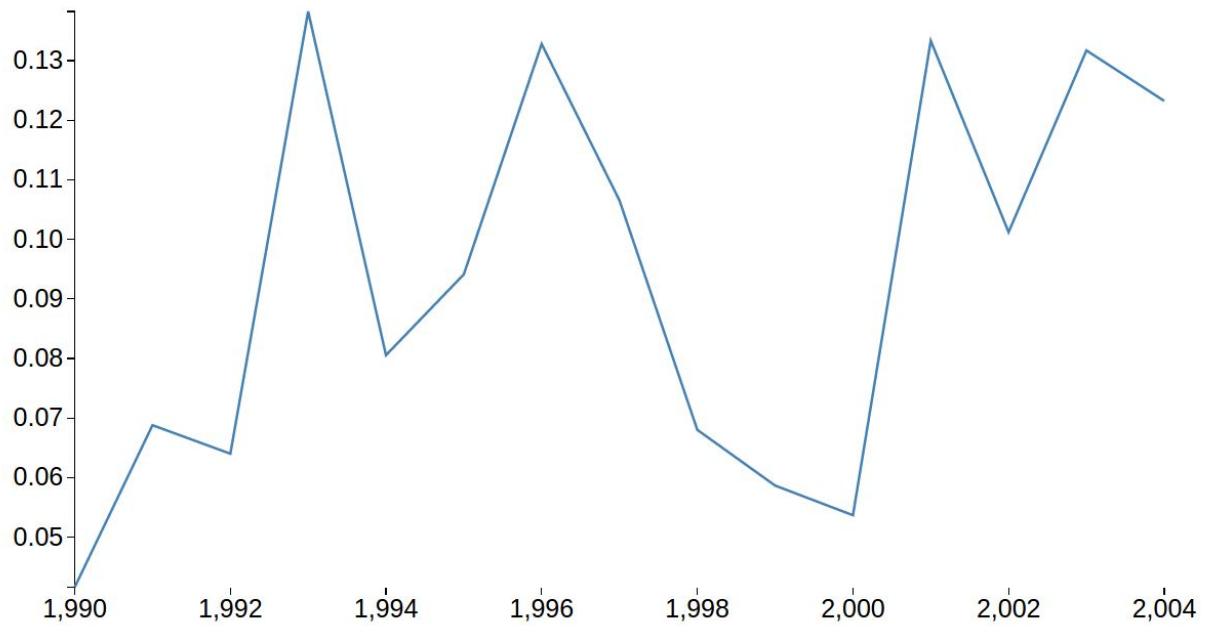


This graph shows the possibility that you can see the word "terrorist" in the new york times "immigration" related articles. Very interestingly, the possibility increases dramatically around 2002. That is when the 911 event happened.



This graph shows the possibility that you can see the word "illegal" in the new york times "immigration" related articles. Even though the graph has some sparks around certain years but the overall graph tend to be more general.





Overall, the clustering tells us mainly two important events, the Kosovo War around 1999 and the Elian Gonzalez custody battle around 2000. Surprisingly, what affects the social media like newspaper the most is not the comparatively recent event like 911 but, instead, some of the older events like kosovo war and cuban custody battle and so on.