

shenghan_project4

Shenghan Yu

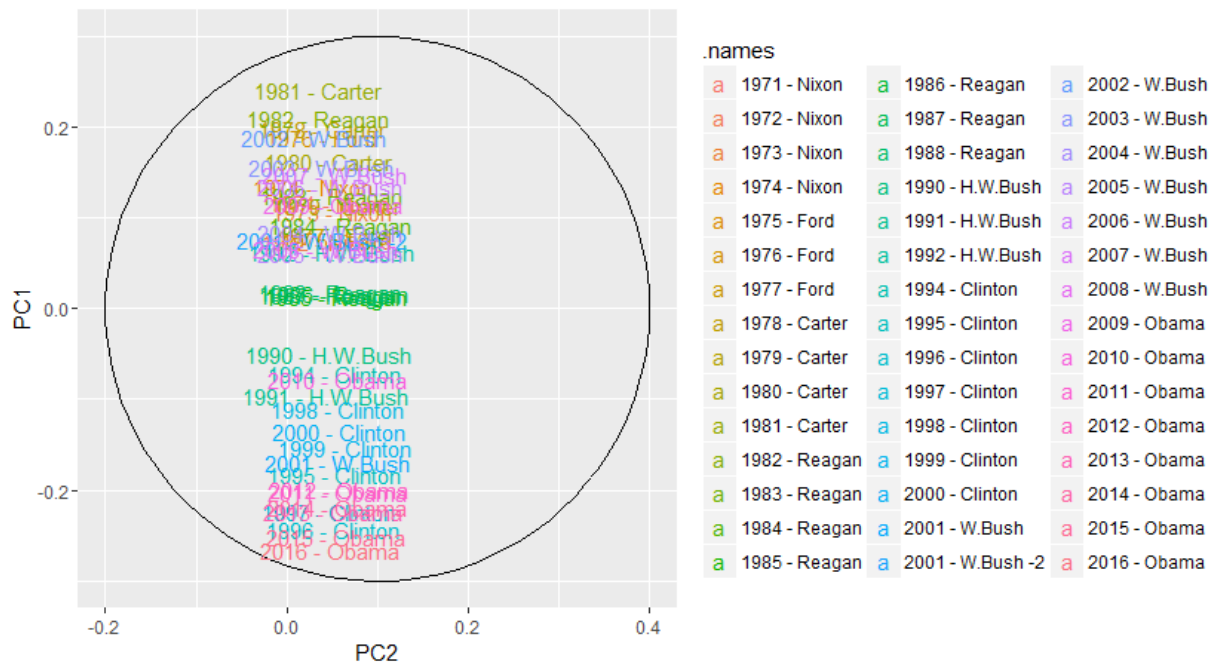
May 12, 2016

PCA Analysis on Corpus

In this section, we conducted a principle component analysis on the corpus of state of union speech. The result shows that the the first principal component explains about 97% of the total variance in the data. We also plotted the all the speeches in the 2D space formed by the first two principal components.

There are several patterns that is worth disucssing here. Overall, there are three clusters in this 2D space, one of which consists of President Ronald Reagan alone. Besides, 6 out of 7 Reagan's speeches are in the cluster. This indicates that the style of Reagon's state of union is significantly different from all other presidents in terms of vocabulary and words picking. Another intriguing finding is that the a majority of Bill Clinton's and Barack Obama's deliveries lies in the same cluster, along with those of George H.W. Bush in 1990 and 1991. This may be attributed to the the three presidents' common emphasis on domestic economy and development during there administration, whether in the golden times of 90s or in recovery from the financial crisis in 2008.

In addition, we can see that the first state of union that George W. Bush delivered (in Feb 2001) lies in the same cluster of Clinton and Obama's, while the latter speeches fall into the opposite cluster. A critical factor here is the 9-11 attach happened on Sep 11th, 2001. George W. Bush delivered an "extra" state-of-union to declare a war on terrorism. This event radically changes the major objective of Bush's admisitration from domestic economy to wars, foreign affairs and securities. That's a sound reason why, after 2001, all the Bush's deliveries lie in the other cluster.



Variance explained by the first 10 principal components

