

Data Science for Industry Project 2

Introduction

Exploratory Data Analysis

The data set consists of 30 State of the Nation Address (SONA) speech transcripts from all the presidents from February 1994 til recently in February 2018.

A hypothesis can be made that sentiments of speeches differ depending on the political season. There are 3 political seasons :

1. Pre-Election
2. Post-Election
3. Normal Term

Below is a summary of all the speeches as per the various presidents categorized by the 3 political seasons :

Period	Presidents	Pre-Election	Post-Election	Normal Term	Total
1994	de-Klerk	1			1
1994-1999	Mandela	1	2	4	7
2000-2008	Mbeki	1	1	8	10
2009	Motlante	1			1
2009-2017	Zuma	1	2	7	10
2018	Ramaphosa			1	1
Total		5	5	20	30

Table 1: The number of SONA per president and arranged by political seasons

From the table above the following remarks can be made :

- de Klerk, Motlante and Ramaphosa all have one speech each which will make it extremely difficult to accurately predict given the far higher number of speeches from their counterpart presidents.
- There are far more speeches done during the normal season which will inherently bias the training data towards that season
- Mbeki and Mandela dominate the number of speeches with 10 apiece. This will also inherently bias the training data towards them.
- Pre-Election speeches are evenly distributed across 5 of the 6 presidents whilst post election speeches are dominated by Mandela and Zuma.

Word Distribution

Below are the most frequently used words of all the 30 presidential speeches rescaled according to their respective political seasons :

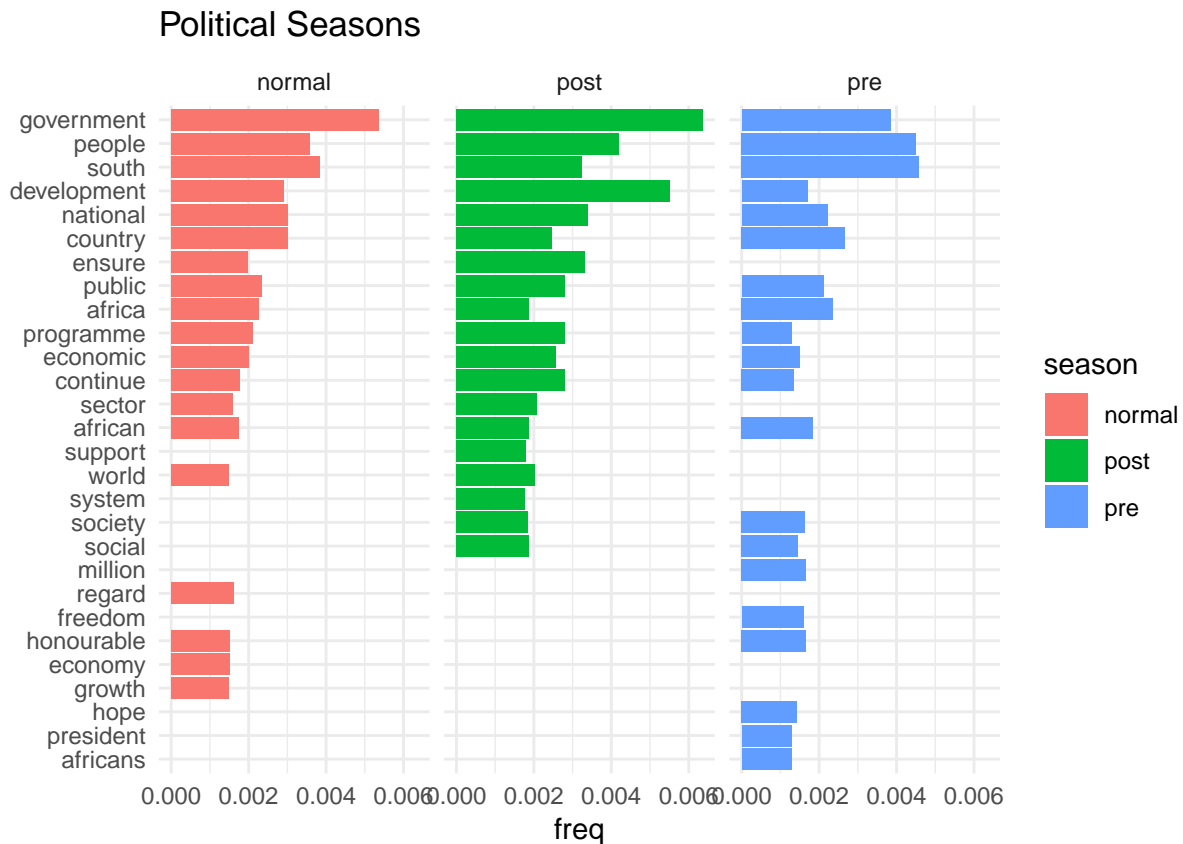


Figure 1: Frequently used words across political seasons

From the illustration above the following remarks can be made :

- Notable words commonly used across all political seasons are **south, africa, government, development, people, country**. These words will potentially not assist in distinguishing the respective presidents.
- Comparing pre and post to normal political seasons, notable words like **economy** and **growth** are introduced into their speeches. The utilization of these words (which could imply economic growth) are understandable given that these are typical themes that need to be addressed constantly throughout the normal period of presidential terms.
- Comparing pre to post political seasons, uplifting words like **freedom, hope, africans** are used before and not after elections.
- Comparing pre to post political seasons, notable words introduced are **support, system, ensure**. These words convey a theme of action and execution which is expected after coming from an election.

Below are the most frequently used words of all the 30 presidential speeches rescaled according to the respective presidents :

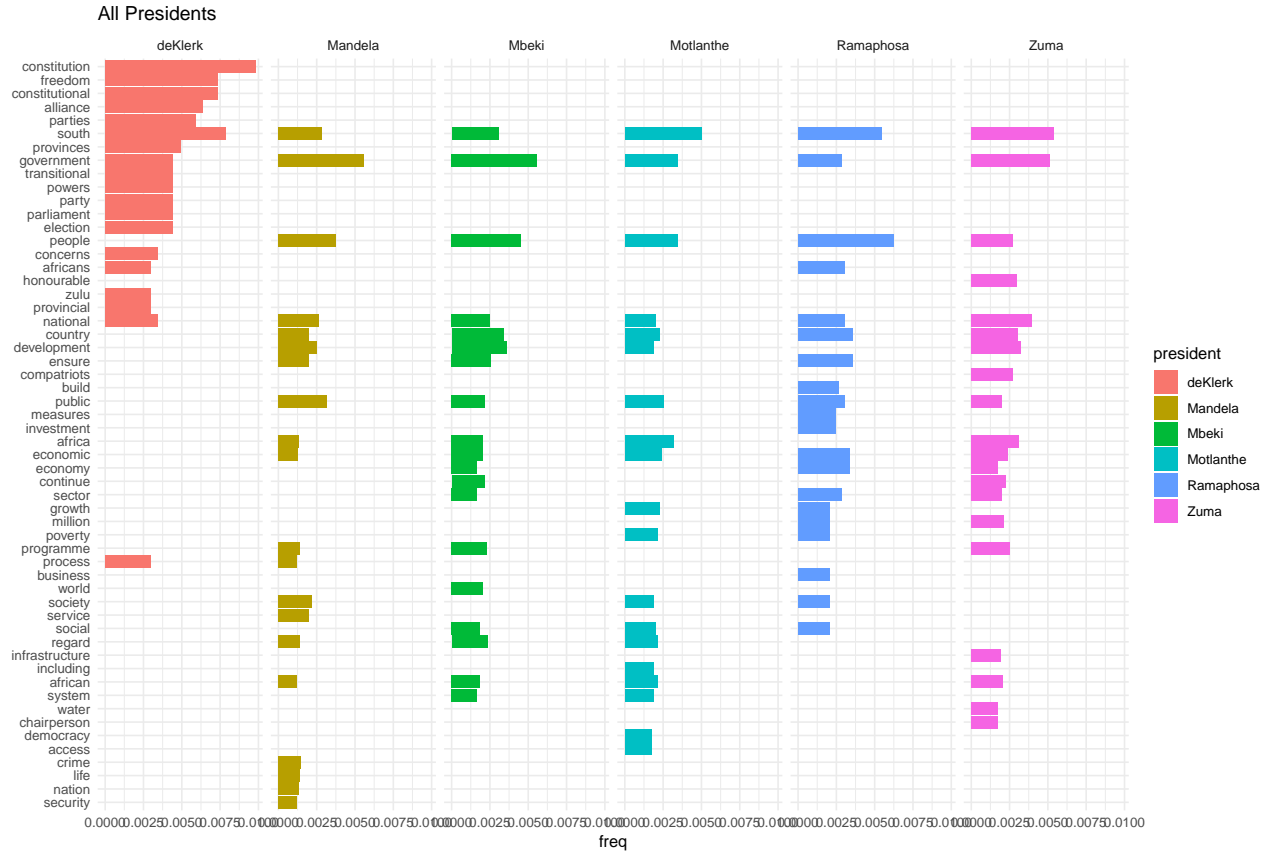


Figure 2: Top words used by all presidents in SONA

From the illustration above the following remarks can be made :

- Looking at all the presidents frequently used words ,de Klerk has the least common words.This is due to the fact that firstly there is only one speech in the dataset and secondly given that the speech was before the first democratic elections,the content will be far different to the speeches made by the presidents that preceded in the democratic era of South Africa.
- Commonly used words across all presidents are **south,government,national** which are also common words across political seasons.
- Notably words commonly used across all presidents excluding deKlerk are **people,country,public,ensure,development**

Clustering by Term Similarity

Words from the respective speeches we aggregated by the respective presidents.Words greater than 4 letters were considered so as to focus primarily on descriptive words.The resultant word counts were then normalized to avoid biases of presidents with more speeches .

K means clustering was then conducted and a $k = 2$ was selected based on the ‘elbow rule’.

The objective is to see what frequent common words do the presidents use and what potential themes to these similar words posses.The resultant visualization can be viewed below:

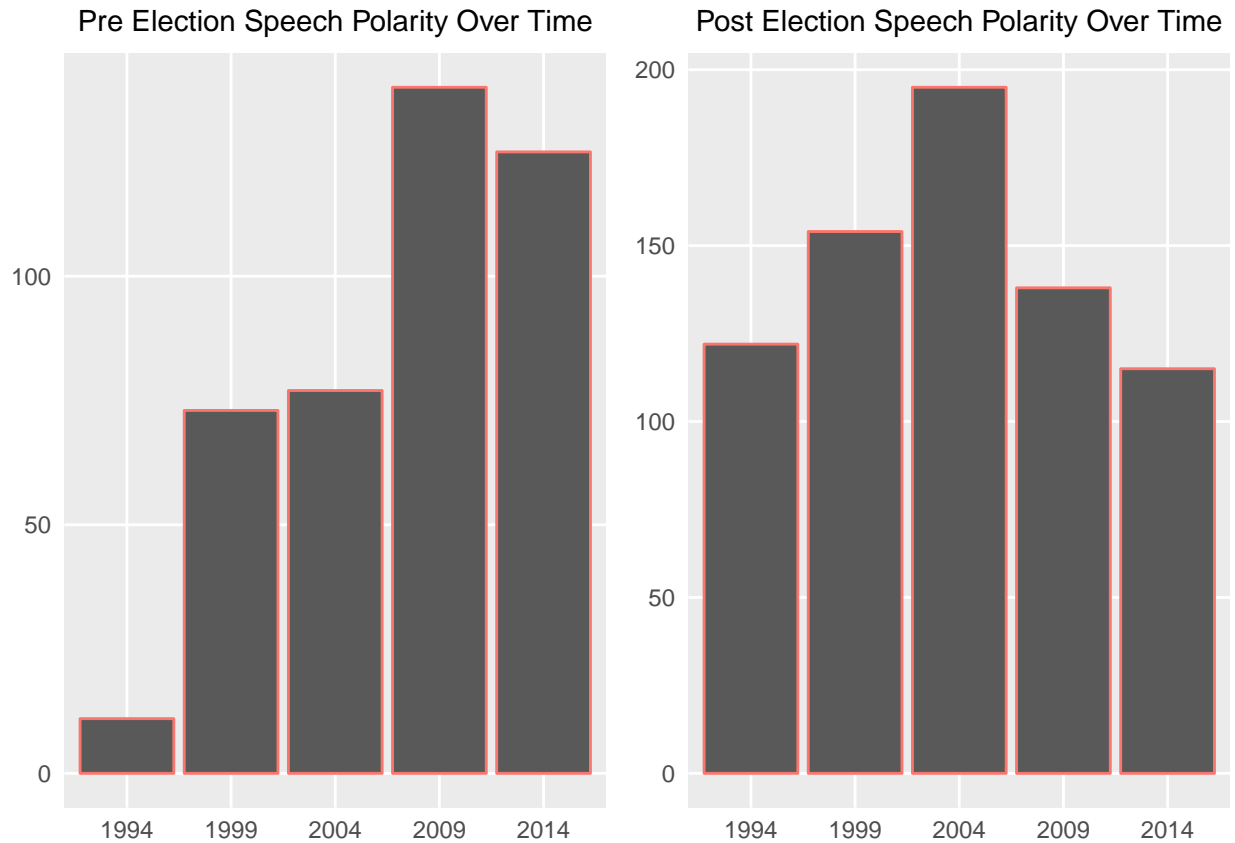
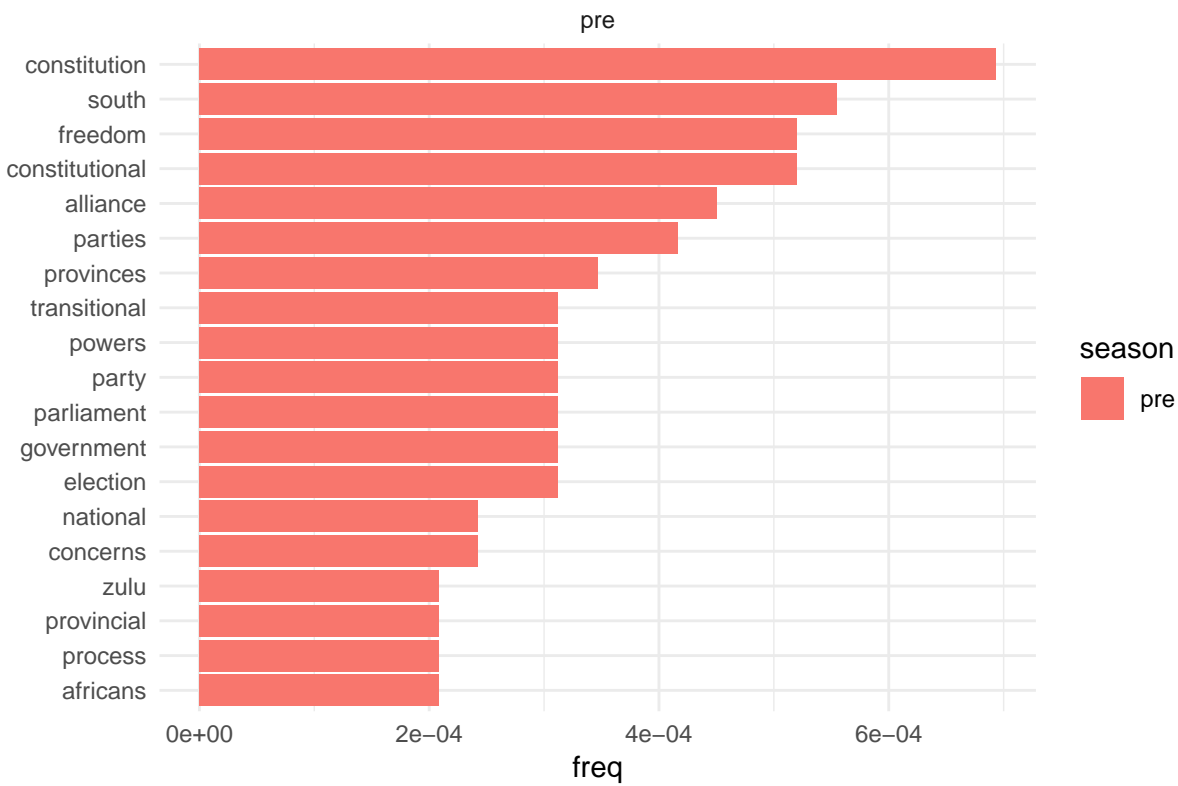


Figure 4: Pre vs Post Election Speech Sentiment over time

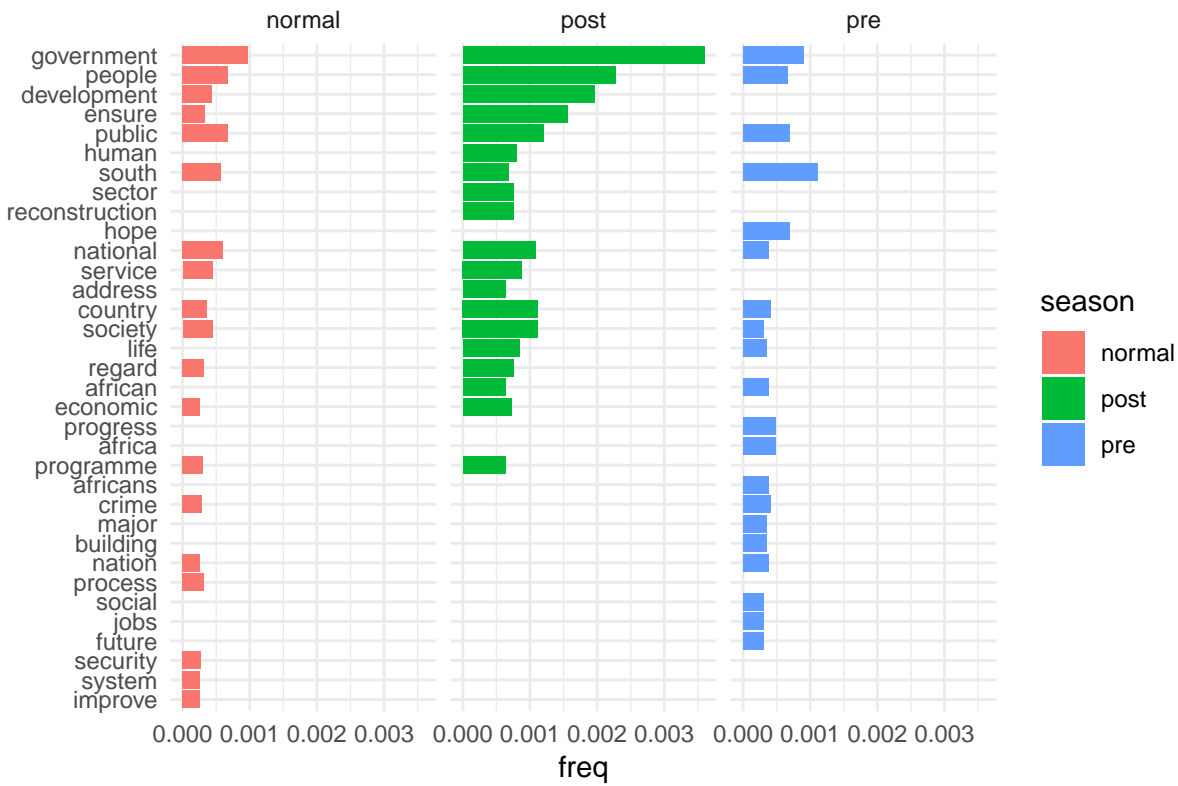
From the illustration above the following remarks can be made :

- There is an interesting behavior in the first 3 election years (1994,1999,2004) there was an improvement in overall sentiment after the elections,however in the more recent 2 election years (2009,2014) there has been a more cautionary tone after elections
- The biggest overall sentiment disparities before and after elections occur in 2004 and 2009.The 2004 pre/post comparison is interesting - both speeches were done by Mbeki and it showed the highest increase in net sentiment in any given year.This could be due to the fact that Mbeki had just won his second election and wanted to reassure the country with a positive speech.The 2009 pre/post comparison is counter intuitive- one would expect that given that when Zuma was elected into power for the first time that he would have a higher net sentiment then that of the his predecessor who was there on a temporary basis.

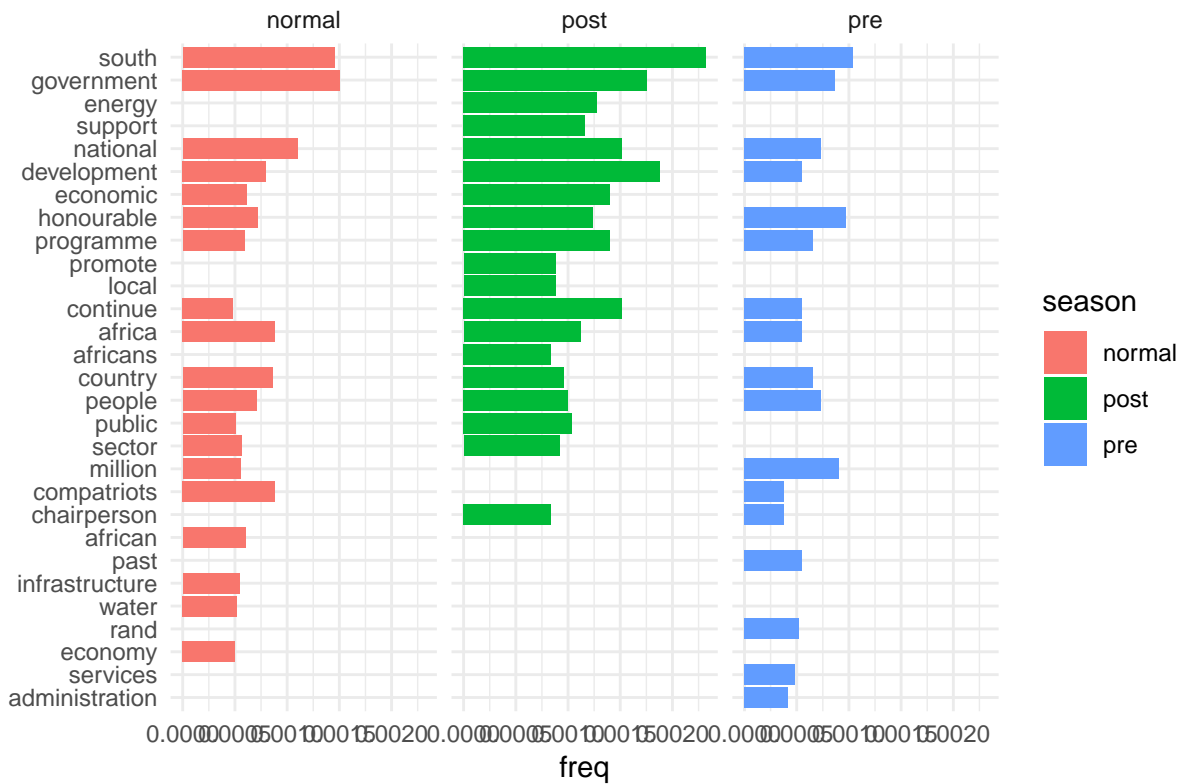
deKlerk



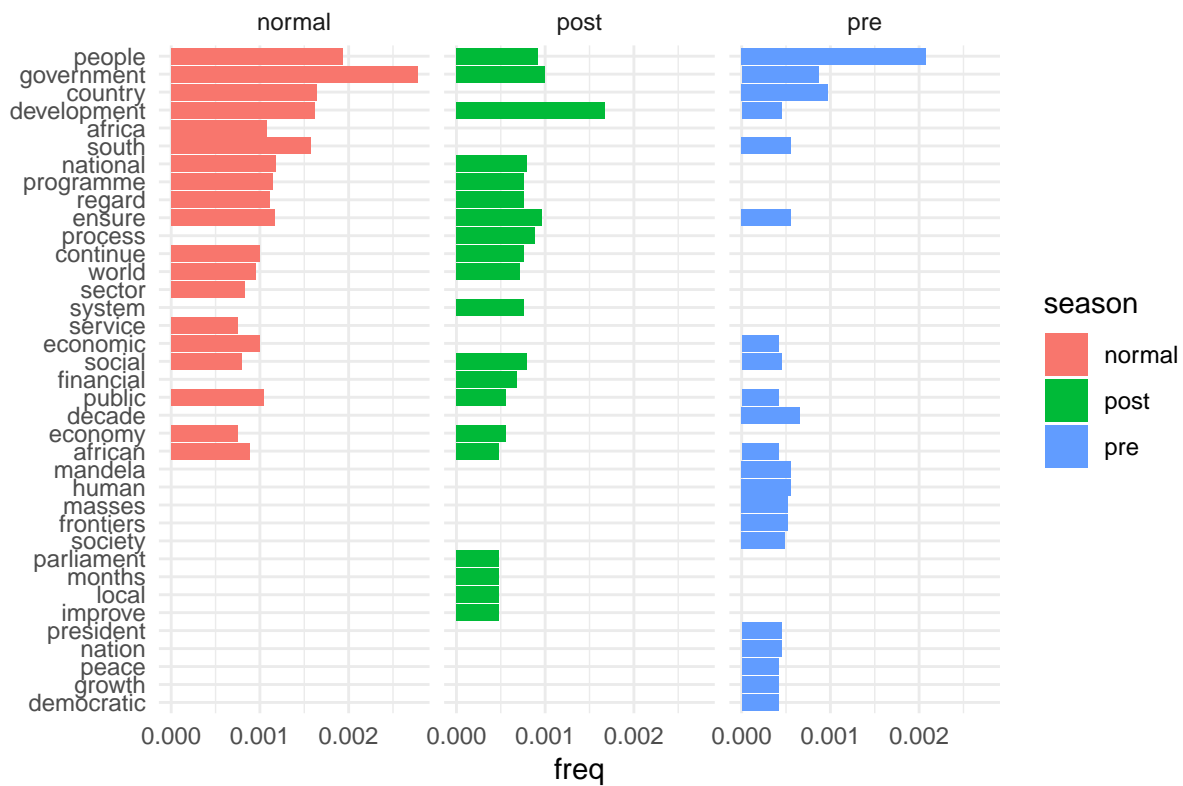
Mandela



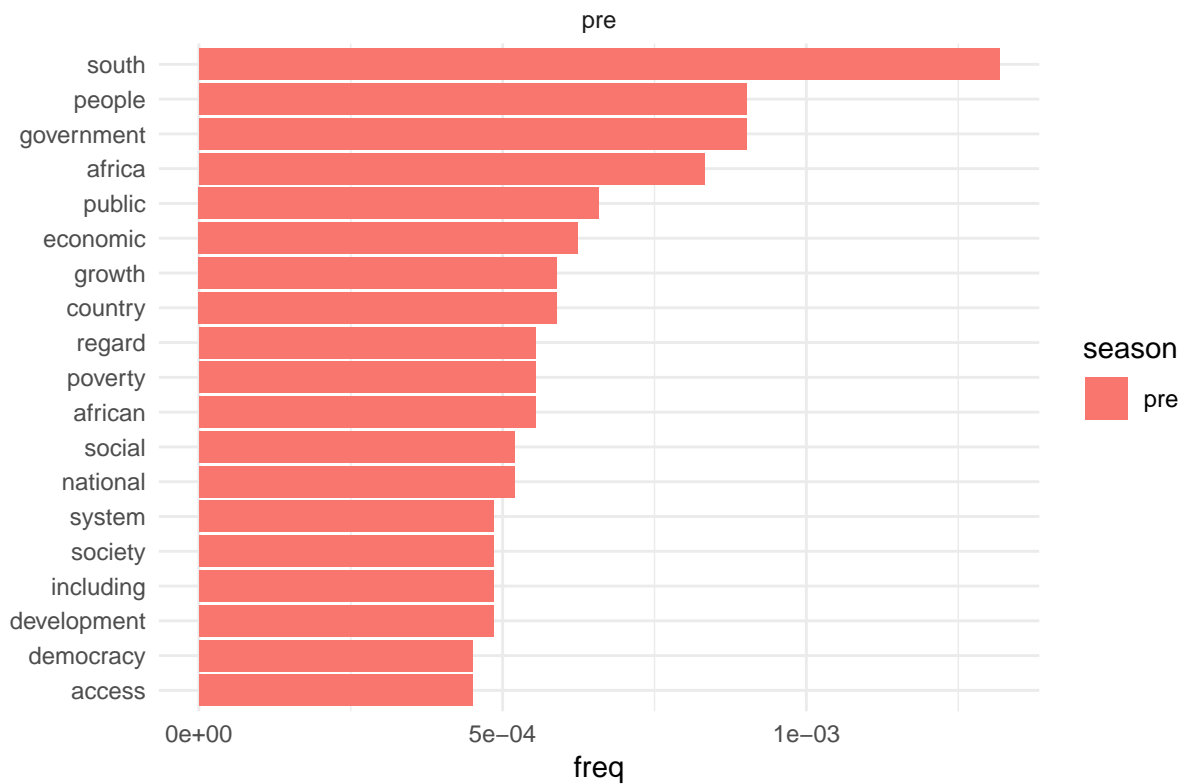
Zuma



Mbeki



Motlanthe



Ramaphosa

