

Topic Modeling + Stance Detection: Kenyan Parliament

Final Report

Sebastian Chacon, Daniel Huff, Eric Long

CS 5322, Natural Language Processing

Professor Lin

May 12th, 2025

Abstract

Contracted through the OIT Research Department at Southern Methodist University, our team has implemented several Natural Language Processing and Machine Learning techniques to connect trends between African Politician Backgrounds to their stances on given topics with politics (ie Finance, Unemployment Rate, Agriculture, Public Education, etc). Our Methodology starts with deriving a set of relationships between our given data being: Biographical Data (web scraped previously by Sebastian), 992 Hansard Text documents from Kenyan Senate, Voting Records (present in ~18% of documents), and 19 given topics from our client. With this information we want to answer the following questions: How does a politician's background affect their decisions on bills in the senate?, Does their association with any given organization affect their voting outcome?, and can a politician's background affect what topic they are for or against?. To answer these questions we followed this methodology: isolate our two way relationship of bills to topics -topics to bills and feed through a Topic Modeling algorithm, then repeat the process except through the one way relationship of speakers (politicians) to stances on topics and feed through a Stance Detection algorithm. After outputting both of these algorithms, we can then integrate the biographical data and deliver a full analysis of trending stances on topics and stances. While this end product was not produced in this report, the project will continue to be worked on and finished post-report.

Introduction

Dr. Cloward in the SMU Department of Political Science is in the process of writing a dissertation on the intersection of international relations and comparative politics, with particular attention to non-governmental organizations, international norms and transnational activism, international aid, Sub-Saharan African politics, and gender and politics. Part of this report requires the completion of a smaller scale project targeted towards the Kenyan Senate and then to be scaled at higher levels through all governments in Africa eventually. Through our testing, we ended with running four models: BERTOPIC for Topic Modeling, base class BERT trained on 2020 Election Tweets for Stance Detection, and a comparison of Logistic Regression and Random Forest for binary stance classification on given pieces of dialogue.

Methodology

This report picks up midway into the research project defined in which we start with the following materials: 992 Documents of Hansard Standardized Texts, 19 Political Topics that we are interested in (ie Agriculture, Education, Unemployment Rate, etc), a list of Speakers from the Kenyan Senate that are notable, and Biographical data on these Speakers (ie College, NCO Associations, etc). As shown in Figure 1, we started by defining our given materials and deriving their relationships. We were able to pull the following relationships: Bills and Topics have a two way relationship in that Bills are associated with many topics and each Topic is associated with many bills. Next we see an equal relationship in that every Text has Speakers that have texts. So they will be treated the same. Finally, a one way relationship of Speakers having stances on any given subject, bill, topic. After establishing these, we broke our project into two approaches: creating outputs for Topic Modeling, and outputting some Stance Detection Models. We took these approaches by implementing SBERT and BERTOPIC for the Topic Modeling, then using a

base class pre-trained (2020 Election Data) BERT and two of our own ML Models: Random Forest, and Logistic Regression.

Starting with Topic Modeling, Sebastian oversaw the implementation of the BERTOPIC model that would take the 992 Texts, perform Optical Character Recognition (OCR going forward) that extracts the text from PDFs, preprocessed our text, and use Semantic Embeddings courtesy of SBERT to use in the BERTOPIC. The preprocessing stage was particularly tough given the input. Starting with the OCR, we had to implement a double layer OCR algorithm due to failure on some levels randomly during runs. Using PLIMAGE it would pick up only portions of texts, scratches on paper would identify as different characters and would ruin our text. We ended up having major success with PyMuPDF that would extract and make all 992 documents usable. After implementing so, there was much time spent on the preprocessing of stop words and stemming. The PorterStemmer would over aggressively stem where words such as “committee” would be stemmed into “comite” which would take away the entire semantic meaning of the word and end up getting removed from our lang detect later through the process. We also had to implement a custom list of stop words for: days of the week, months of the year, politician names, and an 80% threshold for commonly used words. Doing so would allow us to have more accurate and representative topics. From there we used a Lang Detect from NLTK to remove all words that are not in the english dictionary. This serves two purposes, first is to remove any swahili. Although they speak that language often, it is not the official language of Kenya, so we decided to remove it. Second is to remove all misspelled words as they are uncommon enough in unique combinations we can remove this confounding variable as a whole. After preprocessing we feed our clean beautiful data into the SBERT to create topic embeddings and then into the BERTOPIC for topic modeling and labeling.

Before we had started the process for stance detection, we were unaware of any voting records being in the documents. Our plan was to train a model on UK Parliament Hansard documents, which have very accessible voting records. However, after investigating the Kenyan parliament documents, we found that some did have voting records, which is useful for training our models and evaluating their accuracy. Our first method for stance detection was to use a BERT based model specifically trained on stance detection for election data from the 2020 election. At first, this model was evaluating a lot of the text as neutral, which makes sense because parliamentary dialogue is much less concise and polarizing as a tweet can be. We had to force the model to make decisions if it was leaning a certain way to prevent this. After the model predicts each text block as for or against, it counts whether each speaker had more ‘for’ or ‘against’ to predict which way they voted. The model was then tested on the voting records provided at the end of the document.

In our second approach for stance detection, we trained two classifiers, Logistic Regression and Random Forest, on the text of 15 documents from the Senate.csv file that contained voting records. To do this, we first created the necessary training dataset by parsing out the voting records (senator names and whether they voted “Aye” or “No”) from the chosen documents and then labeled each of the spoken texts with how that speaker voted in the end. We then performed standard text preprocessing techniques on these texts, including stop words removal, lemmatization, and the use of a TF-IDF vectorizer to attain TF-IDF embeddings (BERT was not used simply because the length of texts often exceeded the maximum token limit for BERT). We then used this preprocessed data to train the Logistic Regression and Random Forest classifiers, performing 5-fold cross-validation for each model, and scoring them primarily on their Kappa score.

Results

After running the Topic Modeling script, we receive 5 Portions of output shown in Figure 4 a-e. The first output is by document showing what topic it represents (clustered to), the second gives us by topic the representative words and documents and how many documents have been filtered into each topic, the third shows word probabilities by topic for all words embedded, fourth shows us the by document topic probability and boolean for representative document, and our fifth is generated by utilizing OpenAI API to give descriptions of each topic. Looking at Figure 4B we see the BERTOPIC successfully clusters topics together such as: Covid, Pandemic, Hospital, Sickness identifying key concepts from the documents. We can use these in results to connect back to the Stance Detection to create our final desired output of correlation to biographical data. We also see in Figure 4E that OpenAI is able to read and create accurate descriptions by topic of how the documents are connected.

For our first model, we found a very high success rate. The accuracy was typically between 90-100% for any given document that had voting records. An example of a given document's result can be seen in Figure 5. There were some limitations, though. As stated earlier, only 18% of documents had voting records, so we could not evaluate the model's accuracy on the large majority of the documents. In addition to this, the model could only make predictions for those who spoke, not those who voted. This was a major roadblock, as in most documents, the majority of voters do not speak. While this is not a major issue in the larger project, it did prevent some insightful information for our purposes for this portion of the project. Our findings, though, were generally promising for the future of this project.

For the second approach of stance detection, looking at Figure 6, we clearly see that the Logistic Regression model outperformed Random forest in all three metrics; and more importantly, it achieved a Kappa score of 0.546, which suggests moderate agreement between the

Logistic Regression model and the ground truth of the dataset. As such, while there is still much room for improvement in this aspect of stance detection, these results reflect a promising baseline for future work to be built on.

Conclusion

This report showcases moderate success in both the topic modeling and stance detection aspects of the work. While the end product was not fully realized, we have established solid baselines for two major parts of the project. Given that we now have the stances of some speakers, as well as the topics to which the documents belong to, the next step is to connect a speaker's stance with which document they spoke in. In other words, we want to link the document topics to the stances, and by integrating biographical data into this relationship as well, we aim to reveal how a politician's background and affiliations may influence their stance on specific issues, their voting behavior, and the topics they choose to support or oppose.

Figures

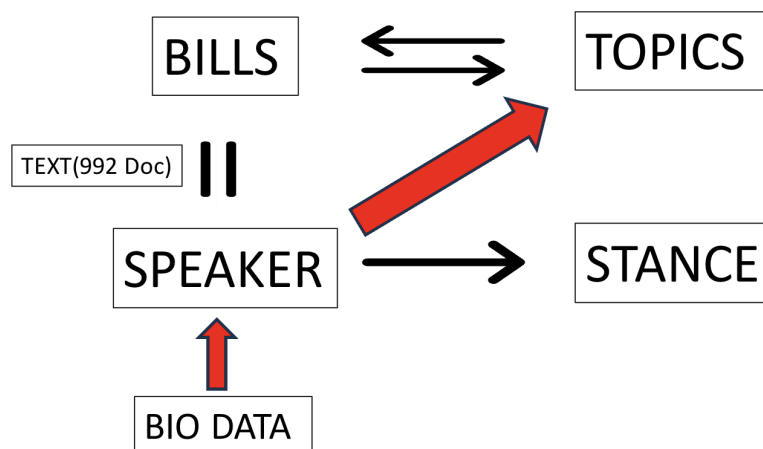


Figure 1:

August 5, 2016	SENATE DEBATES	1
PARLIAMENT OF KENYA		
THE SENATE		
THE HANSARD		
Friday, 5 th August, 2016		
Special Sitting		
<i>(Convened via Kenya Gazette Notice No. 5987 of 2nd August, 2016)</i>		
<i>The House met at the Senate Chamber, Parliament Buildings, at 9.30 a.m.</i>		
<i>[The Speaker (Hon. Ethuro) in the Chair]</i>		
PRAYERS		
COMMUNICATION FROM THE CHAIR		
CONVENING OF SPECIAL SITTING TO CONSIDER THE COUNTY GOVERNMENTS CASH DISBURSEMENT SCHEDULE AND OTHER BUSINESS		
<p>The Speaker (Hon. Ethuro): Honourable Senators, I wish to welcome you back a few days after the commencement of our August recess and thank you for finding time to attend this Special Sitting of the Senate. I am certain that each one of you has a loaded schedule of events at your respective counties. I am sure it is because of the value you attach to the business of the Senate, especially the consideration of key business such as the County Governments Cash Disbursement Schedule, that you have created time to be here to dispose of today.</p> <p>Honourable Senators, the County Governments Cash Disbursement Schedule for the Financial Year 2016/2017 was sent to us by the national Treasury and received at the Senate on Friday, 29th July, 2016, after we had already adjourned for the August recess. Consequently, Hon. Senators, considering the need to approve the Schedule, to enable disbursement of the much needed financial resources to the county governments and the need to conclude consideration of Bills awaiting Divisions at various stages, the Majority Leader, pursuant to Standing Order No.29(1) of the Senate Standing Orders, requested the Speaker to appoint a day for a Special Sitting of the Senate to consider the County Governments Cash Disbursement Schedule for the Financial Year 2016/2017 and other urgent business pending before the House. The request which was made vide Letter Ref.</p>		
<p><small>Disclaimer: The electronic version of the Senate Hansard Report is for information purposes only. A certified version of this Report can be obtained from the Hansard Editor, Senate</small></p>		

Figure 2:

MESSAGE FROM THE NATIONAL ASSEMBLY

REJECTION OF SOME SENATE AMENDMENTS TO BILLS AND APPOINTMENT OF MEDIATION COMMITTEE

The Speaker (Hon. Ethuro): Order Senators! I wish to report to the Senate that pursuant to Standing Order No. 40 (3) and (4), I have received the following message from the Speaker of the National Assembly regarding the decision of the National Assembly on Senate amendments to The Natural Resources (Classes of Transactions Subject to Ratification Bill) (National Assembly Bill No. 54 of 2015) and the Forest Conservation and Management Bill (National Assembly Bill No. 49 of 2015).

Pursuant to the provisions of Standing Order Nos. 41 and 142 of the National Assembly Standing Orders, I hereby convey the following message from the National Assembly:-

WHEREAS, the Natural Resources (Classes of Transactions Subject to Ratification Bill) (National Assembly Bill No. 54 of 2015) was published via the Kenya Gazette Supplement No. 139 of 18th August, 2015 to give effect to the provisions of Article 71 of the Constitution by providing for classes of transactions subject to ratification;

WHEREAS, the Forest Conservation and Management Bill (National Assembly Bill No. 49 of 2015) was published via the Kenya Gazette No. 133 of 11th August, 2015 to give effect to the provisions of Article 69 of the Constitution with regard to conservation and management of forest resources and to repeal the Forest Act of 2005;

WHEREAS, the two Bills were passed by the National Assembly on Wednesday, 9th March, 2016 and Thursday, 17th March, 2016 respectively and referred to the Senate for concurrence;

Figure 3:

Figure 4:

a.

b.

C.

d.

e.

Topic 260:
Keywords: refugee, identity, potatoes, citizenship, camp, somali, insecurity, card, bag, registered
Description: The topic is about the process and challenges of Somali refugees, focusing on their identity and citizenship issues. It may also discuss their life in refugee camps, the importance of having registered identity cards, and basic needs like food (potatoes). The aspect of insecurity might refer to their uncertain and unstable living conditions.

Topic 261:
Keywords: registrar, coalition, expulsion, kang, ata, nomination, ideology, register, indirect, definition
Description: The topic seems to be related to politics or political science. It could involve the process of registering for a political party, the formation of coalitions, the expulsion of members, nominations for positions, and the ideologies or beliefs that guide these actions. It might also discuss indirect methods of influence or the definitions of these political terms.

Topic 262:
Keywords: textile, sugar, cane, clothes, industry, leather, china, company, sucrose, locally
Description: This topic is about the production and manufacturing industries, particularly focusing on textiles, clothing, and leather, as well as the sugar industry which involves the processing of cane into sucrose. It also mentions that these industries are local to China and may involve specific companies.

Figure 5:

```
==== STANCE SUMMARY ====
Total FOR: 208
Total AGAINST: 6
Total NONE: 0

==== SPEAKER STANCES ====
Sen. Wakili Sigeti: FOR
Sen. Mumma: FOR
Sen. Osotsi: FOR
Sen. Maanzo: FOR
Sen. Sifuna: FOR
Sen. Kisang: FOR
Sen. Kathuri: FOR
Sen. Gataya Mo Fire: FOR
Sen. Fakia: FOR
Sen. Khalwale: FOR
Sen. Kinyua: FOR
Sen. Murango: FOR
Sen. Mandago: FOR
Sen. Okenyuri: FOR
Sen. Oburu: FOR
Sen. Methu: FOR
Sen. Veronica Maina: FOR
Sen. Orogeni: FOR
Sen. Wafula: FOR
Sen. Ogola: FOR
Sen. Munyi Mundigi: FOR
Sen. Mwaruma: FOR
Sen. Wamatinga: FOR
Sen. Abdul Haji: FOR
Sen. Tabitha Keroche: FOR

Evaluation Accuracy on voters: 20/20 = 100.00%
Total speakers in vote list: 48
Matched speakers with predictions: 20
Senators who voted but did not speak: {'dullo', 'chesang', 'tabitha keroche to second.', 'githuku', 'munyi', 'kavindu', 'thangw', 'lelegwe', 'methu to second.', 'cheruiyot', 'mungatana', 'madzayo', 'seki', 'sifuna spoke off record', 'lomenen', 'kisang.', 'kavindu muthama', 'mungatana mgh', 'cheptumo', 'ali roba', '}
```

Figure 6:

