

Ayotzinapa in Peril: Gauging narrative shift surrounding 43 “disappeared” students

Anonymous ACL submission

This paper presents an NLP pipeline designed to extract the narratives and analyze the language in a linear timeline surrounding the events of 2014 in Ayotzinapa, Mexico where 43 students were kidnapped. This event is incredibly polemic in Mexico due to the government response and has led to moments of violence and numerous protests, culminating in the firebombing of the Congress building of the State of Guerrero. Tweets dating back to September 2014 up until September 2020 were extracted by the main topics for each month and a timeline surrounding these events was identified; complexity, stylometric, and psychological features including sentiment were extracted from the topics and months identified earlier in the pipeline. Changes in the linguistic features and topics month to month were measured to see if the changes corresponded with an event in the timeline surrounding Ayotzinapa. Using this pipeline, along with the context in which an event happened, a hypothetical analyst can have a better idea of the narrative and linguistic shift surrounding certain events.

Keywords: Ayotzinapa, natural language processing, Web Scraping, sentiment Analysis, feature engineering

1 Introduction

Does Twitter mirror real life, or is it just a stand-alone conduit by which select information is conveyed and has no relationship with real life happenings? Do the narratives propagated through social media (in this case Twitter) display fluctuations in concrete linguistic data while an event happens, or does the linguistic data remain unchanged? The term “narratives” refers to stories or ideas, particularly those of human interest and emotion, that are conveyed in communication. In this case, this study explores what stories and ideas are being propagated regarding the events of Ayotzinapa and

aims to not only extract these stories, but also to see if they correspond to real life and to identify any existing linguistic patterns.

1.1 Background

On September 26, 2014, in the state of Guerrero, Mexico, 43 students from the Ayotzinapa rural teacher’s college were going to a march in Mexico City, but were stopped en route by cartels in collusion with the Mexican government, kidnapped, and “disappeared”. This event was henceforth called the Iguala mass kidnapping, or the events of Ayotzinapa. It became an increasingly controversial topic in Mexico and around the world. It stirred massive protests, violence, and government action that is still felt to this day. This culminated with the most recent 2020 firebombing of the Congress of Guerrero in protest of the government’s involvement. In 2014, the parents of the missing students toured the United States in search of help and to tell the world of their plight; at this event, I personally got to meet these people and join their fight. This study provides another tool by which we can analyze these events through the use of an NLP pipeline and linguistic features and aims to see if the narratives propagated by social media are true to real life events or if they are an inaccurate reflection.

1.2 Relevant Work

There are no exact comparisons to this study in terms of previous work that is relevant to this project. The previous research on Ayotzinapa mostly involves collective trauma and corruption. At the University of Texas at Austin, [Harlow et al. \(2017\)](#) studied the social paradigm shift in terms of media coverage of the events of Ayotzinapa, but linguistic features were not studied and the article focuses more on the engagement of the media. A literature review was completed on linguistic feature engineering, and the work of [Abonizio et al.](#)

(2020) was relevant in choosing what linguistic features might be important in discourse analysis. Lagutina et al. (2019)’s survey of stylistometric features was also referenced. Extensive work has been done on topic modeling and its capabilities; the work of Yang et al. (2014) regarding topic modeling with big data specifically for Twitter, and that of Hong and Davison (2010) were especially helpful in creating the topic model for this study. Ultimately, this is a discourse analysis paper, and the work of Wei-hua (2000) offers a solid foundation that provides insight into the different methods of discourse analysis, including the use of linguistic features and context. This project strives to be an agglomeration of all of this previous work

2 Methods

The question this study aimed to answer was whether one can compare tweets to real life events and extract any meaningful data. Our methodology came to fruition in the form of a pipeline to aid in gathering data, partitioning it into more meaningful information, adding extra features, and finally comparing how the features changed with real life events. Since there was no dataset available related to Ayotzinapa, we needed to collect our own data. We leveraged the SNSscrape package to collect tweets related to Ayotzinapa in order to create a dataset that contains tweets from August 2014 to December 2015. This time period was selected because we wanted to study the events of Ayotzinapa when they were the most relevant and fresh in peoples’ minds while also limiting the amount of data. Next, to extract meaningful data and to lessen the amount of data that we had to study, we used topic modeling to extract the most representative documents for each month. Topic modeling is a common technique used in NLP and has been proven to be accurate and useful in finding topics within tweets and narrowing down information. The topic model used for this study identified the 38 most representative documents/tweets from August 2014 to December 2015. Once these documents were identified, we extracted linguistic features from them; linguistic features are meant to offer quantitative data by which to gauge differences between texts/tweets and how they change over time. This change is ultimately analyzed in conjunction with the events surrounding Ayotzinapa. Discourse analysis is then performed by a linguist or analyst to see if there is any relationship between the



Figure 1: Pipeline

change in the linguistic feature information and the identified event. Our pipeline is shown in figure 1.

3 Pipeline

3.1 Pipeline

The pipeline shown in Figure 1 uses text inputs to produce classification results, isolate prominent topics, and create a final data frame with the outputs.

3.2 Data

The data was collected using SNSscrape and the hashtag “Ayotzinapa” was identified as the primary hashtag to scrape. Other hashtags were identified using <http://best-hashtags.com/> and were found to be “JusticiaParaAyotzinapa” and “AyotzinapaSomosTodos”. We started scraping tweets from August 1, 2014 to get a baseline of activity and features before the events of September 26th occurred in which the students disappeared. The methodology used when scraping was to divide the month into thirds and scrape 100,000 tweets for each third, meaning that a total of 300,000 tweets was collected for each month. We collected data from August 2014 to December 2015, a total of 17 months and 5,100,000 tweets. SNSscrape provides a variety of data that can be used for numerous purposes which are outlined in Table 1. For the purposes of this analysis, we removed any duplicate tweets by doing a lazy match using a pandas function that only kept the first tweet that showed up and dropped the rest. We also dropped any tweets that had less than 10 likes to trim down the data, thus allowing us to keep the most relevant tweets in terms of engagement, propagation, and accordance within the populace.

Table 1 displays the metadata that was scraped using SNSscrape. We used the LikeCount metadata to filter out tweets that did not have more than 10 likes. We used the Content metadata to extract the features and the language used. The Lang metadata was used to distinguish tweets between English and Spanish. That allowed us to only have to create features that handle one language at a time instead of having to use a multilingual model. The Data meta-

| Metadata | | | |
|-----------------|----------------|----------------|----------------|
| Date | User | retweetCount | lang |
| Content | Outlinks | LikeCount | Source |
| renderedContent | tcountLinks | quoteCount | SourceUrl |
| media | retweetedTweet | quotedTweet | mentionedUsers |
| id | replyCount | ConversationID | Source |

Table 1: Snscape metadata

data was used to keep track of when tweets were posted, giving us the opportunity to track change over time.

3.3 Topic Modeling

Following the initial collection of data, we extracted prominent topics from the Twitter dataset. We first preprocessed the data by removing punctuation, words with hashtags, URLs, user mentions, and then removed stop words in Spanish and English as outlined by the NLTK package. This process was completed using a script that uses both the NLTK and Spacy packages. The content of the processed tweets were then lemmatized to produce a corpus for topic modeling. We then fed the corpus through a Latent Dirichlet Allocation (LDA) model where we designated the number of topics we wanted by using a coherent score.

In total, we isolated 38 topics, many of which contained similar keywords (eg. “Ayotzinapa, “mata”), but with different corresponding weights for their respective topics. Subsequent analyses involved isolating topics with atypical keywords and evaluating the most representative text for those topics to gain insight into the narratives surrounding the events of Ayotzinapa for the year. The dataframe at this stage in the pipeline can be seen in figure 2.

| Topic | Topic_Prior | Topic_Coherence | Topics |
|-------|-------------|-----------------|--|
| 0 | 0.0 | 0.5115 | ayotzinapa, congress, cartagena, ay, ya, el, pa... [fista, parry, radison, facidm, , ya, mda, a... |
| 1 | 1.0 | 0.5115 | ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... [interstate, ajerido, facilidad, investigac... |
| 2 | 2.0 | 0.5115 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [f, facidm, fista, ayotzinapa, cub, ayotzinapa, ch... |
| 3 | 3.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 4 | 4.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 5 | 5.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 6 | 6.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 7 | 7.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 8 | 8.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 9 | 9.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 10 | 10.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 11 | 11.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 12 | 12.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 13 | 13.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 14 | 14.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |
| 15 | 15.0 | 0.5026 | ayotzinapa, cub, ayotzinapa, ayotzinapa, ayotzinapa... [fista, ayotzinapa, ayotzinapa, ayotzinapa, ayotzinapa... |

Figure 2: Topic Modeling Dataframe

3.4 Features

Our research focuses on how linguistic features change leading up to certain events. We classify our features into three categories: complexity, stylistic,

and psychological. Abonizio et al. (2020) defines complexity features as linguistic features that capture the overall objective of the context at the word and sentence level. Stylistic features use natural language techniques to gain grammatical information to better understand the syntax and style of the document. Psychological features are closely related to emotions and the cognitive aspect of NLP. In our case we also have an extra feature, sentiment. Sentiment gauges positive or negative emotion towards a person, thing, or feeling.

3.4.1 Stylometric Features

The stylometric features extracted for further analysis were noun counts, adjective counts, adverb counts, and pronoun counts. This was done using the POS function in the Spacy package. The es_core_news_md data was used for Spanish and en_core_news_md was used for English.

3.4.2 Complexity Features

The complexity features extracted for further analysis were character counts, word counts, and type token ratio. The equation used to calculate is show in figure 3.

$$TTR = \frac{V}{N}$$

Figure 3: Type token ratio

3.4.3 Psychological Features

For this paper we did not extract any psychological features.

3.4.4 Spanish Sentiment Analyzer Features

We used Spanish Sentiment Analyzer Python package, which gauges positive and negative sentiment. The closer a text is to 1, the more positive it is, and the closer a text is to 0, the more negative it is. This package uses a sentiment analyzer that was trained using a deep convolutional neural network on 800,000 reviews of users of the pages of eltenedor, decathlon, tripadvisor, filmaffinity and ebay

and presents an 88 % accuracy when determining sentiment.

3.4.5 Final Feature Data Frame

The final dataframe contains 17 sets of data for each month starting in August 2014, from before the 43 students disappeared, all the way to December 2015, providing a glimpse into the narratives for at least a year. They contain the metadata from the initial scrape that can be used in analysis and it also contains the features that have been extracted. We can now see how these features have changed over time, giving us the opportunity to do a qualitative analysis on the features along with a quantitative analysis on how they work in context. Figure 4 shows a concatenated dataframe.

| context | renderedContext | Text | Sentiment | mood_count | work_count | sdi_count | sdi_count | pros_count | cons_count | word_count | word_density |
|---|---|---|-----------|------------|------------|-----------|-----------|------------|------------|------------|--------------|
| ayotzinapa, la punta del iceberg de estados... | ayotzinapa, la punta del iceberg de estados... | ayotzinapa, la punta del iceberg de estados... | 0.497892 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 10.0 | 1.000000 |
| Padres de desaparecidos obligados a revelar... | Padres de desaparecidos obligados a revelar... | Padres de desaparecidos obligados a revelar... | 0.434456 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 118.0 | 9.0 | 11.000000 |
| COPEL, DIRECCIÓN Ayotzinapa no es PALESTE, | COPEL, DIRECCIÓN Ayotzinapa no es PALESTE, | COPEL, DIRECCIÓN Ayotzinapa no es PALESTE, | 0.233842 | 7.0 | 0.0 | 1.0 | 0.0 | 0.0 | 95.0 | 11.0 | 7.916667 |
| Hay 4 refugios por caso de desaparecidos de Ayot... | Hay 4 refugios por caso de desaparecidos de Ayot... | Hay 4 refugios por caso de desaparecidos de Ayot... | 0.389857 | 10.0 | 0.0 | 1.0 | 0.0 | 0.0 | 131.0 | 12.0 | 10.090909 |
| Cronología Caso Ayotzinapa - II Ayotzinapa - II Internacional | Cronología Caso Ayotzinapa - II Ayotzinapa - II Internacional | Cronología Caso Ayotzinapa - II Ayotzinapa - II Internacional | 0.323861 | 15.0 | 0.0 | 1.0 | 0.0 | 0.0 | 108.0 | 17.0 | 9.333333 |

Figure 4: Final concatenated Dataframe

3.5 Features

4 Discourse Analysis

Once the tweets have been collected, the most important tweets and narratives identified, and the features extracted for each tweet, we can now do discourse analysis. To reiterate, discourse analysis is the study of the ways language is used in texts and contexts. We identified prominent events related to the 43 students who have disappeared and analyzed linguistic change corresponding to those identified events.

4.0.1 Timeline



Figure 5: Prominent Timeline of Events

A literature review was conducted to identify the major events that happened in 2014/2015 in regard

to Ayotzinapa. We identified the dates 2014-09-26, 2014-09-26, 2014-10-05, 2014-11-05, 2014-12-08, 2015-01-30, 2015-02-03, 2015-03-15, 2015-04-01, 2015-05-01, 2015-06-01, 2015-07-01, 2015-08-01, 2015-09-01, and 2015-10-01 as important dates. Figure 5 highlights these events in chronological order and Appendix 8.1 at the end of the paper goes into more detail of the events.

4.0.2 Sentiment Observations

Figure 6 highlights the sentiment being at an all time low on September 26, 2014, which is the first date that was highlighted as important in the section above, and on the subsequent dates. This pattern is analyzed for all the important events determined and can be further confirmed in the Appendix.

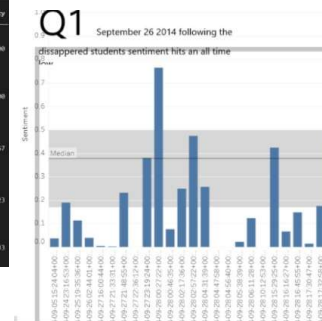


Figure 6: Sentiment to Date Graph

4.0.3 Stylometric and Complexity Observations

In terms of stylometric features in Figure 7, it is apparent that there is a distinct lack of pronouns, which was initially surprising given that people use the word 'I' when expressing solidarity or identifying with the people or the event (ie. I am Ayotzinapa, somos Ayotzinapa). This lack of pronouns was then followed by a sudden, seemingly random, burst in pronoun usage. This phenomenon is attributed to most of my data being in Spanish, but then a tweet in English would be presented as "I am Ayotzinapa". Spanish is a language that drops pronouns, while English tends to keep them. We also analyzed that tweets tend to hover around 20 words, which makes sense in light of the 140 character limit for Twitter at the time. Word density

| Linguistic Features | | | |
|---------------------|-------------|-----------------------|----------------------|
| No | Type | Name | Description |
| 1 | Stylometric | noun _{count} | Noun tag frequency |
| 2 | Stylometric | adj _{count} | ADJ tag frequency |
| 3 | Stylometric | adv _{count} | ADV tag frequency |
| 4 | Stylometric | pro _{count} | PRON tag frequency |
| 5 | Complexity | char _{count} | Number of Characters |
| 6 | Complexity | Type Token Ratio | Lexical diversity |
| 7 | Sentiment | Sentiment | Sentiment |

Table 2: Final Linguistic Features

also seems to be pretty high, which indicates that there is a high lexical diversity and that people find a multitude of diverse ways to announce their support, which mirrors the real life context in which people protested, posted, and rallied around this cause.

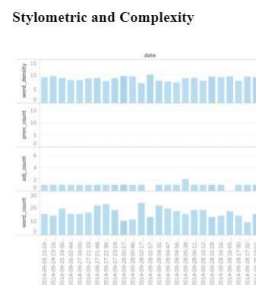


Figure 7: Stylometric and Complexity features to Date graph

4.0.4 Keywords

Keywords remained fairly consistent, except for the final few which differed in accordance with a different event on the list occurring. The keyword Ayotzinapa remained fairly consistent throughout. This means that with the final keywords of the series, we could gauge the shift in narrative. The narrative within the first year rarely shifted, but once one of the events highlighted in section 4.1 occurred, the keywords were related to the text in the timeline, thus highlighting how important those events truly were. This can be seen in Figure ??.

The keywords in September 2014 demonstrate a shift in October 2014 when the next semi-related event happened.

5 Conclusion

In conclusion, we analyzed the language used by the Twitter populace and concluded that the Twittersphere can mirror real life events in terms of

[illegible]

Figure 8: Keyword narrative shift

the change in linguistic features. The stylometric, complexity and sentiment features identified all fluctuated with the real life events in concordance with normal reasoning; narrative shift in terms of keywords and associated content fluctuates similarly, and our small pipeline can help in analyzing such events. This helps by having an abundance of relevant data and features which can be qualitatively analyzed.

6 Future Work

This research provides the foundation which can be used to continue to do an in-depth analysis on Ayotzinapa. The dataset created can be used for discourse analysis, narrative change, and maybe predictive analysis or classification of collective action potential due to the polemic nature of this event and how it has stirred people around the world to action. In terms of future work, I would like to use this dataset to see if I could find a correlation between linguistic features and potential for action. Given enough data, can I predict whether people will protest and to what intensity? I would also like to make my dataset contain more years from 2016 to the present. My limitation with not getting more data was time, CPU/GPU power and just working with that amount of data. My next steps would also be expanding upon the features I already have and creating better features to help in analyses. The final limitation and expansion I would do is create an even better sentiment analyzer for the