

Geopolitical Themed Stressors Extraction from Social Media: A Case Study on Ukraine

Andrea Hannah

Meharry Medical College
School of Applied Computational Sciences
handrea22@email.mmc.edu

ABSTRACT: This case study on Ukraine focuses on emotions and geopolitical related stressors instead of mental health status as some existing studies focus on. This study has demonstrated the potential of using NLP techniques and social media data to explore emotion-based stressors, using EmoRoBERTa to detect and label emotions to each tweet in our dataset and allow the extraction of stressor emotion-based tweets for geopolitical conflicts. We then completed a text analysis involving text cleaning, lemmatization, and TF-IDF tokenization to form a new corpus in which LDA topic modeling could apply dominant topics to a new corpus of cleaned stressor tweet documents, whereby a data science and social science trained professional evaluated topics and performed a subjective and manual extraction of final stressors in the final analysis of the related key terms associated with the assigned topics. Each topic had 30 terms and was reduced to a smaller set deemed sufficient by the professional to conduct the final Stressor extraction that formed the basis of the Stressor and Policy analysis.

Keywords: emotion; emotional health; politics; policy; geopolitical conflict; NLP; Natural Language Processing; Lemmatization; EmoRoBERTa; BERT; Bidirectional Encoder Representations from Transformers; stressors; stressor extraction; term frequency-inverse document frequency; TF-IDF; Latent Dirichlet Allocation; LDA; LDA topic modeling; topic modeling; X; Twitter; tweet; social media; social media engagement; Ukraine; Russia

1 INTRODUCTION

The goal of this research is to identify the sentiments of Twitter Users around the political and humanitarian mood for Ukraine. Russia invaded Ukraine on February 20, 2014 and annexed Crimea, part of Ukraine, in response to the electoral defeat of former pro-Russian President of Ukraine, Victor Yanukovich and the perceived threat to Russia's geopolitical influence in Eastern Europe. This study uses data from the X social media platform from April 8, 2018 as a show case to conduct emotional health analysis and stressor extraction during this disaster event leveraging natural language processing (NLP) and transfer learning techniques. The emotion prediction algorithm employs EmoRoBERTa-based model to analyze emotion behind each tweet, while the stressor extraction mechanism leverages Latent Dirichlet Allocation (LDA) topic modeling for identification of topics, thus, enabling geopolitical stressor themes extraction through the identified topics and the important terms represented in each topic.

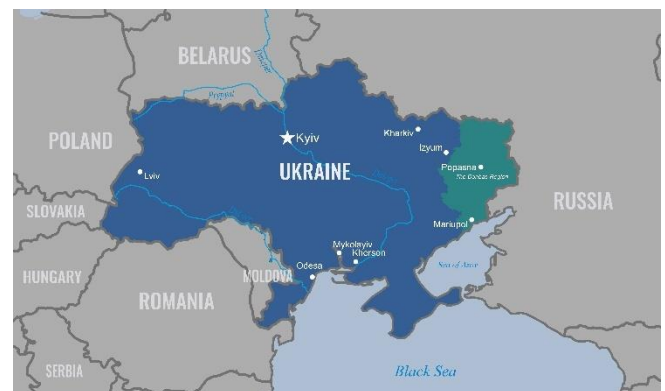


Figure 1

2 RELATED WORK

In this section we review the previous research that inspired many of the NLP techniques used in this work.

In the Bui et al. study, the research team mined emotions and stressors encountered by civilians and shared them on Twitter during Hurricane Harvey in 2017 using a dataset of tweets dated between August 20, 2017 to August 30, 2017 [1]. The dataset consisted of around 400,000 tweets sourced from Kaggle. This study applied a BERT-based model to predict emotions associated with a tweet posted by users before employing natural language processing (NLP) techniques are utilized over negative emotion tweets to explore the trend and prevalence of the topics discussed during the disaster event. Using Latent Dirichlet Allocation (LDA) Topic Modeling, themes were identified which allowed researchers to manually extract stressors termed as climate change related stressors.

The study results showed 20 climate change related stressors were extracted and emotions are peaking at the deadliest phase of the disaster. The study concluded that emotions may be a useful approach for studying environmental well-being outcomes to understand climate change impacts. Similarly, this study for geopolitical stressors should be useful in gaining insight on geopolitical conflicts and sentiments.

2 DATASET REVIEW

The original dataset of 7,390 tweets were sourced from Kaggle and were dated April 8, 2018. The dataset featured tweets relating to Ukraine and Palestine. Having split the tweets by country, the results are file subsets for Ukraine with 3,860 tweets and Palestine with 3,530 tweets. For this study, the focus is on Ukraine for which the data also was reduced again to 834 stressor tweets. The original variables for the dataset include 'ids', 'date', 'user', 'replyCount', 'retweetCount', 'likeCount', 'lang', 'text', and 'country'. In this study, the 'text' column was renamed 'tweet'.

4 PROJECT PIPELINE & METHODOLOGY



Figure 2

To summarize the pipeline in Figure 2, we obtain well-indicated stressors on geopolitical topics, collected tweet data undergoes a text refinement process, where Twitter's emoji, hexadecimal images, special characters, hyperlinks, and unwanted words are being removed.

We start by deploying an emotion classification model to extract tweets with negative emotions to emphasize the stressors' factors. The original 7,390 tweets were split by Country with Ukraine having 3,860 tweets reduced to 834 tweets using EmoRoBERTa and Palestine having 3,530 tweets. Note, the Bidirectional Encoder Representations from Transformers (BERT) structure, namely EmoRoBERTa is utilized to detect emotions. There are 26 emotion labels needing to be tagged to our dataset's tweets; each of the labels represents a distinct emotion for each tweet. After tagging the emotions to the tweets of our dataset, a new csv file is generated with an added variable column, called 'Emotion'. We then work clean this file of 3,860 tweets to extract negative emotions, which will be considered stressors, as we generate a new stressor csv file containing 834 Ukraine related English language tweets. We then move on to Token Featurization and LDA Topic Modeling.

4.1 EmoRoBERTa MODEL

Our first primary objective was to perform emotion detection analysis using EmoRoBERTa to generate emotion labels on a batch sample of our data in order to predict the emotions labels to be applied to the rest of the dataset, using a TensorFlow GPU device processor. This step forms the basis of our analysis and is the source of where our stressors are generated from. After performing the below TensorFlow emotion label prediction function, we will generate a new CSV file without having to repeat this time consumptive step in other parts of our analysis.

```
In [11]: with tf.device('/device:GPU:0'):
          labels = []
          for item in tqdm(batch_tweets, desc="Processed"):
              preds = [lb['label'] for lb in emotion(item)]
              labels += preds

Processed: 100% ██████████ 483/483 [1:08:52<00:00, 8.56s/it]
```

Figure 3

4.2 NATURAL LANGUAGE PROCESSING TEXT ANALYSIS

With the Emotions csv file, we proceed with Natural Language Processing (NLP) for text analysis by starting with text cleaning, using NLTK or the Natural Language Toolkit and clean text libraries to apply functions that remove urls, emojis, punctuation, and double spaces. Then we apply the lemmatization process to avoid redundancy due to words with the same meaning represented in different forms (e.g: go, going, or went are the same verb, with different form due to grammar). Finally, we are able to isolate the positive emotions in our Emotions csv file using sklearn to extract 10 of the 26 Emotion tags, which were labeled as "approval", "desire", "admiration", "love", "gratitude", "excitement", "optimism", "joy", "amusement", and "neutral" from our dataset, leaving what are considered to be stressor emotions and reducing our Emotions dataset from 3,860 tweets to a new corpus of 834 Stressor Emotion labeled tweets filed as a csv. Then in the next step, we updated the English stop words that are commonly used but contributed little to no meaning for the context, such as "a", "an", "the" and so. Additionally, words that were redundant to

the region or prompted place names and people were removed before moving on to tokenization.

4.3 TOKEN FEATURIZATION (TF-IDF)

Token featurization using TF-IDF or the calculation of the term frequency-inverse document frequency (TF-IDF) score for each token was applied instead of word2vec because TF-IDF performs well with larger dataset given that the resulting large corpus would have more terms would have unique weights. Term Frequency measures how often a word appears in a document, and Inverse Document Frequency measures how common or rare a word is across all documents in the collection in a calculation that takes the logarithm of the ratio of the total number of documents in the collection to the number of documents that contain the word plus one. Common words like 'if', 'that', and 'the' that appear in many documents will have a low IDF score, while rare words that appear in only a few of the documents will still have a high IDF score [3].

TF and IDF Formulas:

TF = Occurrence of a word in a document / the number of words in the document

IDF = Log [the total number of documents / (the number of documents containing the word + 1)]

4.3 LDA TOPIC MODELING

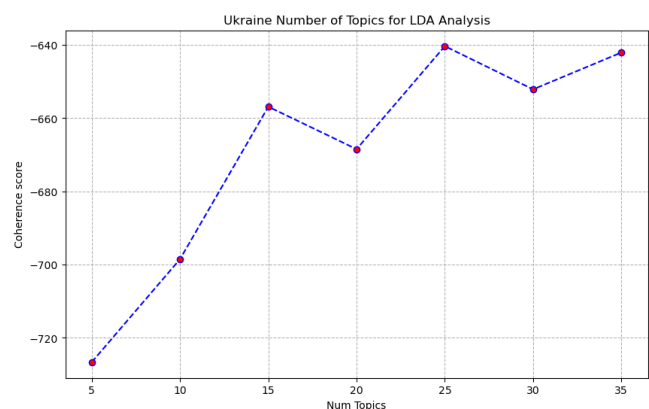


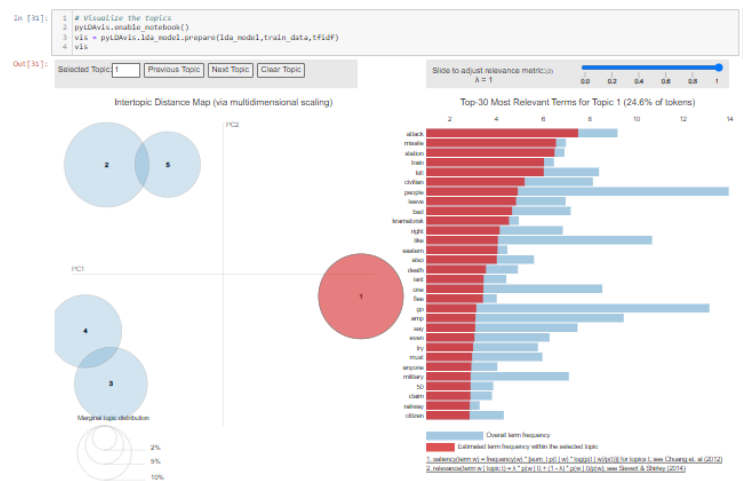
Figure 4

	Topic0	Topic1	Topic2	Topic3	Topic4	dominant_topic
Doc0	0.480000	0.480000	0.590000	0.480000	0.170000	2
Doc1	0.480000	0.480000	0.290000	0.480000	0.540000	4
Doc2	0.480000	0.480000	0.160000	0.480000	0.660000	4

In Table 2, we see the distribution of our 834 stressor tweets among the LDA topics. Topic0 has the highest count of Stressor tweets at 225, second highest count belongs to Topic3 at 207 Stressor tweets, third highest is Topic1 at 149 Stressor tweets, fourth highest is Topic0 with 145 Stressor tweets, and

Topic Num	Num Documents
0	225
1	207
2	149
3	145
4	108

Using NumPy array, we recall the keyword for these 5 Dominant Topics to create a new csv file for Ukraine Stressor Tweets with Dominant Topics. Finally, we visualized the LDA Topics with 30 of its salient terms, and it appears Topic1 has the greatest distance between other topics and straddles two quadrants of the graph, as seen in Figure 5.



We can refine the LDA model to form explicit topics by fine-tuning the model’s hyper-parameters so that we can manually extract stressors on geopolitical topics more precisely.

The second word cloud focuses on stressor emotions, and the 'narrative', 'fear',

‘industry’ as of more importance here. ‘Ruin’ and ‘sacrifice’ along with ‘base’ and ‘Americans’ were hot topics. ‘Inflation’ concerns are raised, ‘morality’ and ‘lies’, ‘peaceful’ and news networks are mentioned, including ‘CNN’, ‘Paul Mason News’ and ‘Disclosetv’. An emphasis on ‘difference’, ‘love’, ‘sell’, ‘loss’, ‘toll’, ‘translation’, and ‘land’ are made. ‘Rescue’, ‘evacuate’, and ‘courage’ reappear in this second word cloud.

5.2 LDA TOPIC MODELING STRESSORS THEMES’ RESULTS, NAMING, RANKING, EVALUATION, AND POLICY RECOMMENDATIONS

Rank by Tweet Count	Stressors Theme	Terms	Tweet Count
1	Topic0: Citizen Death Toll	attack, missile, station, train, kill, civilian, people, leave, bad, death, flee, go, try, must, anyone, military, claim, railway, citizen	225
2	Topic3: Donations and Covid Supply Prices	support, people, god, real, question, please, donate, live, lose, get, must, covid, supply, price, country, fight, need, help	207
3	Topic1: Safety and American Influence	safe, world, life, make, see, help, good, vote, understand, army, time, child, wrong, soldier, come, start, pay, americans	149
4	Topic2: Money and European Influence	know, send, join, invade, money, nato, country, people, buying, keep, tell, think, military, believe, time, anything, fact, government, support, talk, europe	145
5	Topic4: Energy Resources	world, give, gas, oil, stop, next, family, us, innocent, help, long, need, crime, get, heart, support, force, push, energy, start, disgust	108

Table 4

Here are the extractions of the geopolitical stressors in terms that have been reduced from 30 to a variety of numbers, which weren’t reduced to a specified set, such as from 30 to 10. The key term reduction was based on retaining words that added value or meaning and eliminating the ones that do not. The set of key terms were evaluated from a social science perspective to obtain and name 5 Stressor themes to replace our topic numbers. The 5 Stressor themes derived from the LDA Topic Modeling key terms were grouped together along with their former topic number, evaluated, and presented with respective policy recommendation(s) in the order of rank by tweet count per Stressor theme, starting with the highest count being ranked 1 and so on, as shown in Table 4, above. The Stressor theme

analysis and resulting policy recommendations are below:

- ‘Citizen Death Tolls’ (Topic0) is ranked 1 of 5 with a count of 225 Stressor tweets related to this Stressor theme. Given this Stressor having the highest ranking for Stressor tweet counts, there is a very high concern for ‘Citizen Death Tolls’, which indicated substantial and alarming numbers concerning loss of life.
 - Policy recommendation(s): From the Anglophone nations’ perspective, collectively calling for a ceasefire in the interest of preserving human life is necessary and must be repeated and championed. Additionally, a Marshall Plan for Eastern Europe with all but total absolution, will likely become a necessity in concluding this conflict.
- ‘Donations and COVID Supply Prices’ (Topic3) is ranked 2 of 5 with a count of 207 Stressor tweets related to this Stressor theme. Given this Stressor has the second highest ranking for Stressor tweet counts, there is high concern for ‘Donations and COVID Supply Prices’, which indicated need for relief packages, charitable aid, and support for people’s needs. It is not evident why COVID ended up in the terms list because this dataset features tweets from April 8, 2018, before the COVID-19 pandemic. However, since, Ukraine is connected in a large bicontinental landmass of Eurasia then perhaps, COVID was relatively known by these populations with higher contact of populations where COVID is said to have originated. This term seems strange or seems to be an anomaly at first glance, because it doesn’t fit the known timeline of the researcher on when COVID became widespread. However, if the disease reach Russia and/or Ukraine in 2018 then its effects could’ve been just as devastating to vulnerable population, and populations in

wartime conflict is vulnerable. Therefore, it was justifiable to retain this key term for this Stressor theme.

- ‘Safety and American Influence’ (Topic1) is ranked 3 of 5 with a count of 149 Stressor tweets related to this Stressor theme. Here, ‘safe’ is the first key term and includes other terms, such as ‘world’ and ‘vote’, which paired with the appearance of ‘americans’ as a key term, reflects American global influence and support of democratic institutions such as voting. Therefore, we can conclude that this stressor, ranked at 4 of 5 for English language twitter users on Ukraine, have a moderate concern about voting and world input or reaction to this conflict, led by the ‘Americans’ and their influence in global politics and economy.
 - Policy recommendation from the Anglophone nations’ perspective would be to collectively enter talks with the United States and evaluate voting integrity and other democratic institutions in the region to provide consensus on Eastern European Democracy Development plan to the United Nations in consultation with US, British, Canadian, and Australian - Russia policy experts and advisors to draft a plan that is the least threatening to Russia as a pacification strategy until conflict resolution becomes imminent.
- ‘Money and European Influence’ (Topic2), is ranked 4 of 5 with a count of 145 Stressor tweets related to this Stressor theme. Given this Stressor has the second highest ranking for Stressor tweet counts, there is slightly moderate concern for ‘Money and European Influence, which is very reasonable to consider fluctuating exchange rates and how the conflict can trigger inflation since both Ukraine and Russia have different currencies, the hryvnia and the ruble, than the rest of the majority of Europe and European Union member, where the Euro is king. Additionally, since Ukraine is geographically in Europe, considerable concern about how this conflict could spread is very tangible.
 - Policy recommendation(s): The European Union members along with interested Anglophone nations, unless otherwise involved, may consider forming a discovery committee on whether biological weapons have been produced and released to Ukraine and intentionally or inadvertently cause the COVID-19 Pandemic as a result. Unfortunately, the common thread between COVID and Russia is China. Therefore, this is a potential matter of great importance. Even more so at this juncture for the direct conflict at hand, is the quality of the relationship between Ukraine and the EU and Europe, as the EU is a central part of policy, economics, and information sharing on the continent. Relations must not be strained in order for Ukraine to retain or garner more support. For Russia, it must leverage itself to maintain its sphere of influence and gradually de-escalate the conflict in a way that reinvigorates its base around the Kremlin and creatively reinterprets its position with improved moral reasoning behind its actions regarding Ukraine. Claiming Ukraine as an ancestral homeland is not accepted by Anglophone nations, the EU, led by Germany as is supported in our Word Cloud in Figure 7. The concept of kinship or ancestral similarity may be severely diminished or eroded as effective propaganda with the local Russian constituencies should the conflict persist, and the ruble remain severely devalued, as is the hryvnia. Russia is in a unique position to transform Eastern Europe if it rapidly de-

escalates armed conflict and negotiates a bid with the EU and America with a Marshall-like plan for the 21st century, effectively setting the stage to re-enter Global Politics as a partner rather than an aggressor if it can end the conflict peacefully as supported in the Stressor word cloud in Figure 7.

- ‘Energy Resources’ (Topic4) is ranked 5 of 5 with a count of 108 stressor tweets related to this Stressor theme. Given this Stressor’s low ranking for Stressor tweet counts, there is less concern for ‘Energy Resources’ than any of our other Stressor themes. However, given we intentionally reduced the number of stressors to 5 for time efficiency in this study, the theme should still be considered of great importance, especially given, that Ukraine is home to nuclear energy plants and has experienced nuclear energy disasters before on April 26, 1986, where the entire town of Chernobyl had to be abandoned due to its radioactivity. Similar disasters could strike again if armed conflict results in a nuclear power plant being significantly damaged or hit by missiles, artillery, etc. Additionally, gas and oil are listed in the key terms for this Stressor, and this is to be expected as gas and oil prices fluctuate in a region where it is sourced and produced for global and local consumers, especially given bans and trade restrictions that have been places on Russia over time during this conflict.
 - Policy recommendations(s): Engage the scientific community on methods to satisfactorily decontaminate Chernobyl and other radioactive disasters should they arise with staff and robots.

6 LIMITATIONS AND CONCLUSIONS

Besides time, computational power, and current X research access limitations, here are certain limitations of the study. This 2014 data was self-reported by Twitter’s users which may introduce social desirability biases. The focus on Twitter’s data solely may not represent the emotional health of individuals across all web platforms or in real life, limiting the generalizability of our findings. The extraction of stressors manually through main themes discussed on social media may have human bias and dependencies on the identified themes which may neglect stressors in "low volume" discussion. Finally, our model stems from the absence of context information for each tweet. Unlike other social media platforms, tweets from 2018 can only contain up to 280 characters, and a substantial amount of context is concealed within URL links and emojis in shared messages. This can potentially lead to the misclassification of emotions associated with each tweet, thus limiting our capacity to identify more profound stressors. However, the X platform now has larger available characters, so part of this limitation can be resolved by using current datasets once access limitations are addressed and opting to retain emojis in the NLP processing.

There are three technical improvements to consider for this study, which include retaining emojis in the Natural Language Processing steps; increasing the feature count in the tokenizer steps of TF-ID, since there are more tweets in the dataset than the default feature number of 500 represents; and increasing the Topic Number from 5 to a higher optimal number in the LDA Topic Modeling steps.

For future work, we will extend the algorithm to more thoroughly to other geopolitical conflict impacting people with more current tweets, extending the study to Palestine, which hopefully, the findings in such as study could provide insights for policy makers to handle the needs that arise for real-time geopolitical conflicts where related stressors can be reviewed in the policy recommendation process for stakeholders.

REFERENCES

- [1] Bui, T., Hannah, A., Madria, S., Nabaweesi, R., Levin, E., Wilson, M., Nguyen, L. (2023, December). Emotional Health and Climate Change Related Stressors Extraction from Social Media: A Case Study in Hurricane Harvey. Mathematics. <https://www.mdpi.com/journal/mathematics/>.
- [2] Oumaima, I. (2021). Palestine-Ukraine. Kaggle. Retrieved September, 2023, from <https://www.kaggle.com/datasets/oumaimaidhika/palestine-ukraine/>.
- [3] Zhu, Aaron. "Understanding Tf-IDF and Cosine Similarity for Recommendation Engine." Medium, Geek Culture, 4 Apr. 2023, <https://medium.com/geekculture/understanding-tf-idf-and-cosine-similarity-for-recommendation-engine-64d8b51aa9f9>.