

Analyzing the Impact of Domestic Violence on Social Media using Natural Language Processing

Krishna More
Xavier Institute of Engineering
Department of Information Technology
Mumbai University
krishnaaaxo@gmail.com

Frason Francis
Xavier Institute of Engineering
Department of Information Technology
Mumbai University
frason.kalapurackal@gmail.com

Abstract— Due to the rapid advancement in social media and technology, it generates a large amount of data in different areas of applications. Social media analysis and text mining are all about collecting the most valuable data and drawing actionable conclusions. Text mining also referred to as data mining it is which contains various nodes in the form of data which is often linked together to form a pattern. High-quality information is typically derived through the devising of patterns and trends through means such as statistical pattern learning. In this study we have analyzed and mounted social media data from Twitter, new articles, and Reddit which suggest that domestic abuse is acting as an opportunistic infection, flourishing in the condition created by the pandemic. The computing tweet sentiments of domestic violence amongst various social media platforms is a major factor of concern. We have used several topic modeling techniques such as Latent Semantic Analysis (LSA) uses a bag of words model, Hierarchical Dirichlet Process (HDP) is a nonparametric Bayesian model for clustering problems, and Latent Dirichlet Allocation (LDA) is a generative probabilistic model for collections of discrete data. Therefore, in this project, we tend to propose a deeper insight into the rise in domestic violence on social media and to provide a holistic approach to tackle this situation.

Keywords—Sentiment Analysis; Domestic Violence, Opinion Mining; Social Media; Natural Language Processing; Domestic Abuse; Text Mining

I. INTRODUCTION

Domestic Violence is not a pandemic, it's an epidemic. With Covid-19 ravaging the economy; increasing unemployment such crises are set to become much more frequent. In France, reports of domestic violence have increased by 30%. In Spain, the emergency number for domestic violence received 18 per cent more calls in the first two weeks of lockdown than in the same period a month earlier. The United Nations Secretary-General, Antonio Guterres called for urgent action to combat the worldwide surge in domestic violence. He also say that, "I urge all governments to put women's safety first as they respond to the pandemic". This is because domestic violence is an ongoing serious social problem worldwide (Xue et al 2019)[6]. It is estimated that one-third of women worldwide have experienced some form of domestic violence by their intimate partner in their lifetime (WHO 2017).

As this problem has serious precautions on the physical and mental health of the victims (children, woman, man), a lot of research has been already conducted to analyse the nature of

this social problem from interviews with victims, and official data like health records. Recently there has been significant effort put into social media data research. However, much of this research has found many limitations to create practical tools to deal with domestic violence. Within this project, we are looking utilizing topic modeling to assist with domestic violence data classification.

II. LITERATURE REVIEW

A. Topic Modeling

Topic Modeling is a modeling method that aims to find hidden topics in the text. For example, there are two sentences. The first sentence is "Tmall shopping carnival is empty alleys", and the second sentence is "Alibaba's income may hit a new high this year." If we lack prior knowledge and don't understand what "Tmall" and "Ali" are, we cannot find out the relationship between them. However, through TM, we can find that their semantics are similar.

$$\frac{p(\text{word}|\text{speech}|\text{text file})}{\sum \text{the } p(\text{word phrase}|\text{main theme})} * \text{the } p(\text{main topic}|\text{text file})$$

here are two main methods for TM training inference, one is pLSA (Probabilistic Latent Semantic Analysis) and the other is LDA (Latent Dirichlet Allocation). pLSA mainly uses EM (Expectation Maximization) algorithm [11]; LDA uses Gibbs sampling method. A good topic model should produce human-interpretable topics that are distinct from each other. The top words in the topic should coherently belong to some unified concept without too much overlap with the concepts in other topics [8]. These standards are chosen so the "best" topic model is the one that corresponds to distinct subjects of discussion in the particular dataset of concern [8].

There are several different topic modeling techniques:

Latent Semantic Index (LSI), also known as Latent Semantic Analysis (LSA), is a concept put forward in the field of information retrieval. It is mainly to solve two types of problems. One type is polysemous. For example, the word "bank" can refer to banks or river banks; the other type is polysemous, that is, the problem of synonyms, such as "car" and "automobile". The same meaning, if in the search process, when calculating the similarity of the two types of problems, relying on the cosine similarity method will not be able to deal with such problems well. Therefore, a method of latent

semantic indexing is proposed, which uses SVD dimensionality reduction method to map terms and text to a new space. The LSA could also be utilized for information retrieval, which is first implemented in 1988. Since we can get the embedding of documents, we can match documents by calculating the cosine similarity between the documents.

The Latent Dirichlet Allocation (LDA) technique is a generative probabilistic model in which each document is assumed to have a combination of topics similar to a probabilistic Latent Semantic Indexing model. In LDA selection the number of topics is directly proportional to the size of the data, while the number of topic terms is not directly proportional to the size of the data. In Latent Dirichlet Allocation model for text classification purposes the hyperparameter alpha represents density of topics generated within documents and beta represents density of terms generated within topic [12].

Hierarchical Dirichlet process (HDP) is a powerful mixed-membership model for the unsupervised analysis of grouped data. Unlike its finite counterpart, latent Dirichlet allocation, the HDP topic model infers the number of topics from the data [10]. Each set of data is represented by a mixture, with the number of components left open-ended and inferred by the model automatically. Furthermore, components can be shared across groups, allowing for effective modelling of cross-group interdependence as well as generalisation of new models [9].

B. Our Purpose

The online data generation has seen a rapid upward in scale. The noisy and unstructured characteristic of such data has further added to an overall complexity of processing and utilizing the available required information. The mix of messages that are of personal nature in a form of mere opinions or empathetic thoughts along with general awareness promotions have greatly diminished the DVCS services efficacy. With topic modeling, we are looking to identify the different topics or classes of the tweets or comments in reddit platform to make the large and unstructured data more organized in a way that will make it easier for NGOs, government officials or researchers to assess and get useful insights from to better analyze this crisis and come up with better courses of actions accordingly.

At present, there is no available benchmark dataset for domestic violence with pre-labeled data for multiclass classification, so it is a very manual intensive job. Topic modeling could help identify the different topics so it is easier to build a corpus later with labeled data that can serve as a training dataset for deep learning algorithms like CNNs and RNNs for classifying domestic violence texts.

C. Scope of Application

The findings of the twitter and reddit data can be compared to official reports on the prevalence and severity of DV. Official reports and official news concerning domestic abuse can be compared to the findings. Investigation of Demographics: The metadata linked with postings taken from online forums can be a useful source of demographic information on users. Social workers or victim advocates can use our tool when safety planning with victims. He or she can

gauge the level of risk and recommend or intervene appropriately. Secondly, district court judges can either accept the result of the model as evidence to support the victim's case in protective order hearings, or a he or she can use the tool to put a more quantitative assessment of whether a respondent is likely or not likely to commit family violence in the future. Our website itself can be extended to serve as a data gathering tool where survivors can submit anonymous information relevant to their situation. Also, deep learning can be utilized to discover risk factors. And finally, natural language processing can be implemented to provide a chatbot service for people who desire to have a more interactive experience when consuming information but are not yet ready and hesitant to talk to a real person.

III. TOLL DESIGN

This is the architecture for our approach. Our platform will be like a repository that contains the results and insights from our sentiment analysis and topic modeling techniques to help develop better research around domestic violence. Also, this platform can provide great insights for NGOs and government officials to help them with their policy making and strategies to help prevent and detect and fight against domestic violence.

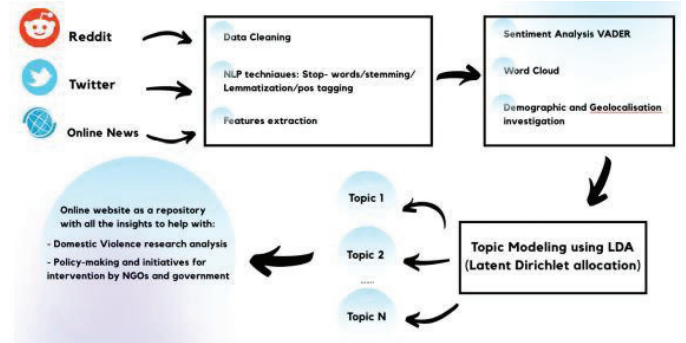


Fig 1. Schematic flow of the project

IV. MODELING METHODOLOGY

A. Methodology

Data is collected and scraped online on twitter, reddit and from online news articles. Specific hashtags used for every file for twitter analyzing the tweets.

- 1 = METOO
- 2 = WHYISTAYED WHYILEFT
- 3 = HeForSheAtHome, WomenCount, GenerationEquility, AntiDomesticViolenceDuringEpidemic, Mask-19 WithHer, SpotlightEndViolence,
- 4 = staysafe, domesticviolence, DomesticAbuse, DomesticViolence

B. Packages Used

To perform sentiment analysis, we used VADER (Valence Aware Dictionary and sentiment Reasoner) it is a sentiment analysis tool that is specifically designed to examine sentiments expresses in social media. VADER uses a

combination of A sentiment lexicon is a list of lexical features (e.g., words) which are generally labelled according to their semantic orientation as either positive or negative [13].

We decided to use VADER since it has been found to be quite successful when dealing with short pieces of text on social media and product reviews. This is because VADER not only talks about the Positivity and Negativity score but also tells us about how positive or negative a sentiment is.

We used Genism and pyLDA visualization packages in order to perform our topic modeling and pyLDAvis helped with creating the visuals of the clusters of topics we found. Those visuals are saved as HTML files and attached with this report.

C. Evaluation Metrics

Perplexity along with topic coherence measure was used for the purpose of evaluation. Where perplexity measures indicate how “surprised” the algorithm is to see a term within a given topic (lower values indicating a better model) and is measured as the normalized log-likelihood of a held-out test set (Mishra, 2018)[1]. However, numerous authors have suggested that it is an ineffective metric because it does not capture semantic information. For our analysis it will be calculated, but the output will not be used to evaluate the model.

Topic coherence measures score a single topic by measuring the degree of semantic similarity between high scoring words in the topic (Mishra, 2018)[1].

This is a measure composed of four steps namely, segmentation of words subsets, probability estimation, confirmation measure then aggregation. Moreover, the measure it is calculated between 0 and 1, where 1 is the most coherent.

V. RESULTS

A. Reddit Data Analysis

Through all the comments, mainly women are seeking help, sharing their abuse stories, how they survived and how they want to help other people. One of the main insights is that the most important topic is about family abuse, not relationship abuse, mentioning mother, father, relationship and child as well as kid, we assume that these people are young adults or teenagers living with their families. Time is also a concept that is frequent among all topics, quoting “year”, “month”, “day” as important words describing the testimonies.

- Recommendations
 - Organizations and NGO should have initiatives and communication through Reddit since it is a space in which people share their testimonies, is a channel where the initiatives can be communicated to the right people.
 - Partner abuse is highly related to partner abuse.

B. Twitter Data Analysis

Through all the tweets the main conversation is from organizations and movements to invite people to join the

movements through hashtags and facts, we can also perceive not only organizations but sponsors.

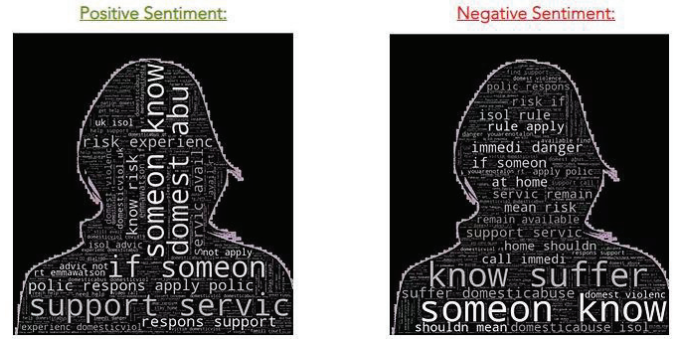


Fig 2. Shows Analysis of hashtags: #staysafe, #domesticviolence and #DomesticAbuse.

Women that tweet about these hashtags and are marked as positive, using “Vader Sentiment Library” in Python, the main topic is women seeking support about their past experiences, another highlight of the words “someone known” in both positive and negative wordings.

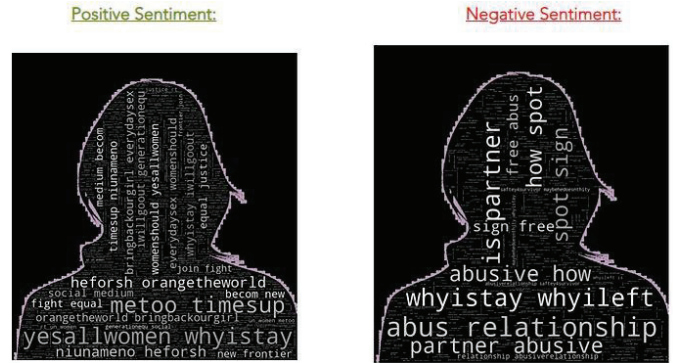


Fig 3. Shows analysis of hashtags: #WHYLEFT, #WHYISTAYED

With the analyzed hashtags, which were part of an initiative that has been present since 2014, stand out movements that have been related to the topic: #metoo, #heforshe, #yesallwomen #timesup, #iwillgoout, #bringackourgirls, and the most recently #niunamenos, led by México in March, also the organization “orange the world” stands out which is a campaign from UNESCO fighting violence against women.

TABLE I. SHOWS THE RELEVANT INFERRED TOPICS:

Topics	Hashtag used	Relevance
Topic 1	#WHYLEFT, #METOO, #WITHHER	This movement revolves around raising awareness and acting with the organization - IPU parliament (global organization of national parliaments to promote peace, democracy and sustainability), also the sponsorship of women in sport.
Topic 2	#METOO, #believewoman	Victims sharing their stories and support of one another and sharing their survivor stories, organization - IPU parliament (global organization of national parliaments to promote peace, democracy and sustainability), also the sponsorship of women in sport
Topic 3	#genderequality	Feminism and equality - Backed by the

		organization "everyday equality" to drive equal social opportunities in UK.
Topic 4	#timesup	Campaign to report and end violence thus creating an equal future, promoting also job quality and activism.
Topic 5	#oranjeworld #yesallwomen	Abuse regarding rape in family and work – organization like Orange the World: Generation Equality Stands against Rape.
Topic 6	#womenshould #everydaysexism #niunamenos	Promotes discussion of psychological abuse, harassment and solution making.

- Recommendation
 - An improvement that can be done in twitter is not only promote de initiatives and join the conversation but invite victims to share their story via Reddit, because there is no character limit, thus increasing Reddit's testimonies and having more information regarding this subject.

C. Twitter Data of top 5 countries

After identifying the top 5 countries that are most engaging in the conversation about domestic violence on twitter, we wanted to further identify and analyze the topics in each of those countries using topic modeling.

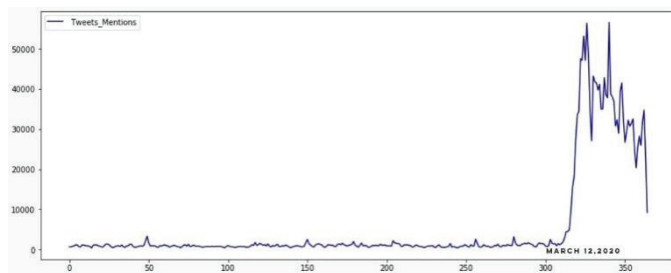


Fig. 4 By analysing the hashtags: #staysafe, #domesticviolence, #DomesticAbuse, #whyistayed, #whyileft, #metoo, #heforsheathome, #womencount, #antidomesticviolenceduringpidemic, mask-19 the trend has elevated since the lockdown.

1) India:

The clustering of topics when filtering only on India are diverse:

- Topic 1: The first topic is about gender equality and the laws against domestic violence. there is a mention of Prasad who is an Indian lawyer who recently wrote about male domestic violence act. This topic is basically about having sometimes biased laws against domestic violence.
- Topic 2: This topic is about assistance and help and society support.
- Topic 3: The 3rd topic is about the pandemic and lockdown and how it is increasing the problem of domestic violence.
- Topic 4: This relates to organizations, initiates, and government actions like UNICEF, NC India (National Committee for Woman), PMO India (prime

minister of India who created a care fund for emergency and distress situations)

- Topic 5: Relates to topic 1 in the sense that it also focuses on “make home safe for man” as well as “woman”

2) United Kingdom:

UK topic are not only about organizations, but the topics revolve around mainly domestic abuse, trying to aid women, in the second topic inviting women to call immediately during COVID.

- UK topic are not only about organizations, but the topics revolve around mainly domestic abuse, trying to aid women, in the second topic inviting women to call immediately during COVID.
- Covid19 help: Aiding women to call during the actual crisis and trying to aid women.
- Charities.
- Community support and sharing stories of victims.
- Aid during COVID promoting a helpline and communicating that support service remains available.
- Same as topic 1 and 2.

3) United State of America

US Is the News country, where is more about the actual sharing of news and the story of the movement. With Joe Biden recent abuse scandal, the first 5 topics revolve around it, also there is a mention of “Alyssa Milano” American Actress that started the “metoo movement” and is in support of Biden.

The only “Corona Lockdown mention” is “Covid” and #Stayathome, but there is no communication towards women or stories about women in any of the hashtags, the only identified movement is #metoo and #believing_woman and other related to Joe Biden and Senator Tara.

4) Kenya:

Topics clustered from filtering on Kenya are oriented towards gender-based violence and gender equality, and also initiates related to DV:

- This topic is about gender equality in the case of domestic violence.
- 2nd topic is about UNWoman organization and support to domestic violence victims.
- This topic is about the domestic abuse and education.
- This is about initiatives related to DV like #withher that talks about creating a world where woman and girls are safe, and “eachforequal” movement that was created in international woman day 2020.
- same as topic 6 it is about GBV(gender based violence) by the United Nations Population Fund (UNFPA) which is something Kenya has recently been shining a light on.

5) *Nigeria*:

Topics identified in Nigeria are as follows:

- This topic is about movements like IDefendHer_heforshe and NGOs in Nigeria like Ceen foundation with the mission of promoting public safety, security and accessible justice.
- The 2nd topic is about community sharing and support for victims.
- COVID19 pandemic and lockdown.
- This topic is about child marriage about domestic abuse.
- Gender equality and UNWoman.
- Organizations and movements to provide aid and support to woman who suffer from Domestic Violence during this pandemic like Global Spotlight that support women's organization that are working to prevent and address violence against women.

VI. LIMITATIONS

One of the challenges with news articles is that unless it's opinionated, there is little room for a diverse identification of topics. But our results generated topics which had great separation.

The noteworthy topics generated are:

Topic 1: Children's services: words most prevalent in the topic were "victim", "child", "resource", "home" and "risk". By identifying the percentage contribution of the most dominant document, main stories are about the lack of these services.

Topic 2: Survivor stories: words most prevalent in the topic were "woman", "household", "victim", "survivor" and "impact".

Topic 3: Police response: words most prevalent in the topic were "city", "police", "victim", "family" and "shelter".

As started previously, there were more topics generated but clearly identifying with the limited number of articles generated (only 26 articles), created a barrier to adequately class the topics. For future work, one can enrich the data with more news articles. For visualizations and individual insight, refer to our website,

<https://ffkalapurackal.wixsite.com/stopdomestiviolence>

VII. CONCLUSION

By raising awareness, sharing knowledge, and bringing voices and data to the public, social media has become increasingly important in the fight against violence [3]. This study conducted can be analyzes the discussion of domestic violence on social media platform. Topic modeling and sentimental analysis was used to extract different subject of discussion, sentiment analysis was used to classify tweet as positive or negative. In topic modeling several text mining techniques like Latent Semantic Analysis, Hierarchical Dirichlet Process and Latent Dirichlet Allocation were used to identify the spikes in topic discussion which directly corresponds to real world events. Finally, the words analysis from the domain-specific embeddings enable us to obtain

greater insight into the abuse types as well as the conditions experienced by the victims [3].

VIII. FUTURE WORK

Based on the conclusions, the following recommendations for future work were made:

- Explore the use of Guided LDA to pre-seed the data words which frequently occur.
- Increase data collected from online news articles and to further enrich corpus to possibly help with clearly identifying the topics.
- Diversify the data sources and media platform and improve the preprocessing techniques applied across the data collections.
- Deal with the bias of creating our dataset possibly including more random sample data.
- Incorporate temporal topic modeling to see changes of the topics over time.

REFERENCES

- [1] Mishra, A. (2018) Towards data science, Metrics to Evaluate your Machine Learning Algorithm. Available at: <https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-fl0ba6e38234>
- [2] Nathan, A. J. and Scobell, A. (2012) 'How China sees America', Foreign Affairs, 91(5). doi: 10.1017/CBO9781107415324.004.
- [3] Subramani, S. et al. (2019) 'Deep Learning for Multi-Class Identification from Domestic Violence Online Posts', IEEE Access, 7, pp. 46210–46224. doi: 10.1109/ACCESS.2019.2908827.
- [4] Subramani, S. and O'Connor, M. (2018) 'Extracting Actionable Knowledge from Domestic Violence Discourses on Social Media', ICST Transactions on Scalable Information Systems, 5(17), p. 154807. doi: 10.4108/eai.29-5-2018.154807.
- [5] Subramani, S., Vu, H. Q. and Wang, H. (2018) 'Intent Classification Using Feature Sets for Domestic Violence Discourse on Social Media', Proceedings - 2017 4th Asia-Pacific World Congress on Computer Science and Engineering, APWC on CSE 2017, pp. 129–136. doi: 10.1109/APWCConCSE.2017.00030.
- [6] Xue, J., Chen, J. and Gelles, R. (2019) 'Using Data Mining Techniques to Examine Domestic Violence Topics on Twitter', Violence and Gender, 6(2), pp. 105–114. doi: 10.1089/vio.2017.0066
- [7] 30 Questions to test a data scientist on Natural Language Processing [Solution: Skilltest – NLP], Available at: <https://www.analyticsvidhya.com/blog/2017/07/30-questions-test-data-scientist-natural-language-processing-solution-skilltest-nlp/>
- [8] Dahal, B., Kumar, S.A.P. & Li, Z. Topic modeling and sentiment analysis of global climate change tweets. Soc. Netw. Anal. Min. 9, 24 (2019). <https://doi.org/10.1007/s13278-019-0568-8>
- [9] Hierarchical Dirichlet Processes, Author(s): Yee Whye Teh, Michael I. Jordan, Matthew J. Beal and David M. Blei, Source: Journal of the American Statistical Association, Vol. 101, No. 476 (Dec., 2006), pp. 1566–158, Available at: https://www.stat.berkeley.edu/~aldous/206-Exch/Papers/hierarchical_dirichlet.pdf
- [10] Hierarchical Dirichlet Process Genism documentation Available:<https://radimrehurek.com/gensim/models/hdpmodel.html>
- [11] Automated Duplicate Bug Report Detection Using Multi-Factor Analysis, Jie ZOU, Ling XU, Mengning YANG, Xiaohong ZHANG, Jun ZENG, Sachio HIROKAWA DOI: <https://doi.org/10.1587/transinf.2016EDP7052>
- [12] Principled Selection of Hyperparameters in the Latent Dirichlet Allocation Model, Journal of Machine Learning Research 18 (2018). Clint P. George, Hani Doss; 18(162):1–38, 2018.
- [13] Simplifying Sentiment Analysis using VADER in Python (on Social Media Text) by Parul Pandey (2018) Available at: <https://medium.com/analytics-vidhya/simplifying-social-media-sentiment-analysis-using-vader-in-python-f9e6ec6fc52f>