

REPORT ON
SOCIAL UNREST PREDICTION USING
SOCIAL MEDIA

by
NIKHIL BIRADAR
20025561

CONTENT

- Abstract.
- Introduction.
- Background & Literature Review.
- Methodology
- Emotion Classification
- Results
- Citations.

ABSTRACT

This study employed natural language processing techniques to analyze sentiments and emotions expressed in online discussions related to social movements, specifically focusing on subreddits like BlackLivesMatter, stopasianhate, and immigration. The investigation utilized two distinct approaches: sentiment analysis leveraging NLTK Vader and BERT, and emotion classification using LSTM-based deep learning models.

The sentiment analysis revealed a nuanced distribution of sentiments within these discussions, with positive, negative, and neutral sentiments expressed across diverse topics. Further, the BERT-based sentiment analysis demonstrated an overall positive trend within the sentiments expressed in the analyzed comments.

Moreover, the emotion classification model effectively identified and categorized emotions conveyed in textual data. The LSTM model trained on a labeled dataset accurately predicted emotions, displaying the potential of deep learning in understanding complex emotional expressions in text.

In addition to sentiment and emotion analysis, topic modeling using LDA was employed to extract key themes across discussions. This revealed various topics prevalent in these conversations, providing deeper insights into the underlying issues and concerns driving these online interactions.

The study underscores the significance of computational approaches in discerning sentiments, emotions, and topics within online discourse related to societal movements. These methodologies illuminate the multifaceted nature of online discussions, shedding light on the diverse emotional expressions and themes prevalent in social movement discourse on Reddit.

INTRODUCTION

From the Black Lives Matter movement to protests against anti-Asian hate and discussions surrounding immigration, online forums like Reddit have become focal points for societal discourse and activism. These digital spaces serve as arenas where individuals express sentiments, emotions, and thoughts, shaping conversations around crucial societal issues.

The proliferation of social movements on digital platforms echoes a global trend of using online spaces for social and political engagement. This paradigm shift from traditional forums to digital platforms has enabled individuals to voice opinions, share experiences, and participate in discussions spanning a wide spectrum of societal concerns. The nature of these discussions, encapsulating emotions, sentiments, and diverse viewpoints, underscores the need for computational tools to comprehend and analyze these multifaceted conversations.

This study delves into sentiment analysis, emotion classification, and topic modeling of discussions within subreddits associated with social movements. Leveraging Natural Language Processing (NLP) techniques, sentiment analysis tools like NLTK Vader and BERT, as well as deep learning models for emotion classification, were employed to dissect the underlying sentiments and emotions embedded in textual data.

The study also involved an LSTM (Long Short-Term Memory) model implementation utilizing GloVe (Global Vectors for Word Representation) embeddings. This model was designed to classify emotional tones within textual data, dissecting nuanced emotions expressed in social media discussions.

This paper contributes to the understanding of computational approaches in discerning sentiments, emotions, and topics within online discussions. It underscores the significance of employing these tools to decipher the intricate

nature of societal discourse and highlights the nuanced expressions and themes prevalent in discussions related to societal movements on platforms like Reddit.

The subsequent sections detail the methodologies used, the findings derived from sentiment and emotion analysis, topic modeling insights, and the significance of computational techniques in dissecting online discussions related to societal movements.

BACKGROUND & LITERATURE REVIEW

Role of Social Media in Societal Discourse:

The advent of social media platforms has revolutionized communication, providing avenues for individuals worldwide to express opinions, engage in discussions, and catalyze social movements. Platforms like Reddit have become pivotal in shaping contemporary discourse, serving as hubs for discussions on societal issues, politics, and activism. Research indicates the significant role these platforms play in reflecting and influencing public sentiment and behavior (Gruzd et al., 2018).

Sentiment Analysis and Emotion Classification:

Sentiment analysis techniques, such as Natural Language Processing (NLP) and emotion classification models, have gained traction in comprehending the underlying sentiments within textual data. These methodologies aid in quantifying emotional tones expressed in online discussions, allowing for nuanced interpretations of user sentiments towards various topics (Pang & Lee, 2008).

Topic Modeling and its Significance:

Topic modeling techniques, such as Latent Dirichlet Allocation (LDA), provide a systematic approach to extract prevalent themes or topics within textual datasets. These models aid in identifying underlying patterns in discussions, unveiling the latent structure of discourse, and categorizing content based on prevalent themes (Blei et al., 2003).

METHODOLOGY

❖ Experiment Design:

The study employs Python-based tools and libraries for data collection, processing, and analysis. PRAW (Python Reddit API Wrapper) is utilized for data retrieval from specified subreddits related to social issues like 'BlackLivesMatter,' 'StopAsianHate,' and 'Immigration.' Sentiment analysis leverages the NLTK Vader Lexicon and BERT-based sentiment analysis using Hugging Face's Transformers. Emotion classification utilizes machine-learning approaches based on text data.

❖ Analysis of Reddit Comments for Sentiment and Topic Modeling:

The analysis focuses on extracting sentiments and identifying topics from Reddit comments related to specific subreddits ('BlackLivesMatter', 'stopasianhate', 'immigration'). This report outlines the methodologies employed to conduct sentiment analysis and topic modeling on the comments retrieved from these subreddits. The methodology involves retrieving comments, performing sentiment analysis using NLTK Vader and BERT, visualizing sentiment distributions, and employing LDA-based topic modeling. Results showcase sentiment percentages, sentiment distribution charts, and identified topics, contributing insights into community sentiments and thematic discussions within these subreddits.

❖ Deep Learning Model for Sentiment Analysis:

➤ Data Preprocessing:

The code loads a dataset ('Emotion_classify_Data.csv') containing comments and their associated emotions. The 'Comment' column is converted to lowercase, and a text-cleaning function is applied to remove URLs, special characters, punctuation, digits, and stopwords using regular expressions and NLTK's word tokenization.

➤ Encoding and Tokenization:

The emotions are encoded using Scikit-learn's LabelEncoder, converting categorical emotions into numerical values. The text sequences are tokenized using Keras's Tokenizer, which converts each text to a sequence of integers.

➤ Word Embeddings (GloVe):

The code utilizes pre-trained GloVe (Global Vectors for Word Representation) embeddings to represent words as dense vectors in a high-dimensional space. GloVe captures semantic relationships between words by mapping them to vectors in such a way that similar words appear closer in the vector space.

➤ Model Architecture:

The neural network model is structured using Keras's Sequential API. It consists of an Embedding layer, an LSTM (Long Short-Term Memory) layer, and a Dense layer. The Embedding layer transforms the integer-encoded sequences into dense vectors, initializing with pre-trained GloVe word embeddings. LSTM, a type of recurrent neural network (RNN), processes sequences by maintaining long-term dependencies, capturing information over extended sequences efficiently. The Dense layer with a softmax activation function performs multi-class classification for predicting emotions.

➤ Model Compilation, Training, and Evaluation:

The model is compiled using the sparse categorical cross-entropy loss function and the Adam optimizer. Training occurs over five epochs with a batch size of 64, utilizing the training data while validating on the validation set. Finally, the model's accuracy is evaluated on the test dataset to gauge its performance.

➤ Mathematical Aspects and Formulas:

VADER Sentiment Analysis: VADER computes a compound score based on word polarity and intensity. The compound score formula involves a summation of the normalized scores of each word in the text. VADER (Valence Aware

Dictionary and sEntiment Reasoner) computes a compound score based on word polarity and intensity. The compound score is derived from normalized scores of individual words in the text. The formula to compute the compound score involves a weighted summation of the valence scores of words, accounting for the punctuation shifts, capitalization, and degree modifiers. The compound score

C is calculated as follows:

$$C = \frac{\sum_{i=1}^n S_i \times W_i}{\sum_{i=1}^n W_i}$$

- LSTM: The LSTM layer involves mathematical computations of gates (input, forget, output), activation functions (sigmoid, tanh), and memory cells. The gates control the flow of information, and the memory cells store and regulate data flow across sequences, utilizing complex matrix operations.

GloVe (Global Vectors for Word Representation) embeddings use co-occurrence statistics from a corpus to create word-word co-occurrence matrices. These matrices are then factorized to obtain word embeddings in vector space. The mathematical concept involves minimizing the difference between the dot product of word vectors and the logarithm of the words' co-occurrence probabilities:

$$J = \sum_{i,j=1}^V f(P_{ij}) \cdot (\mathbf{w}_i^T \cdot \tilde{\mathbf{w}}_j + b_i + \tilde{b}_j - \log(P_{ij}))^2$$

- Emotion Classification:

Utilizing machine-learning models for emotion classification, comments were categorized into different emotional tones. The analysis showed prevalent emotions within discussions, including 'anger' (25%), 'fear' (20%), 'joy' (30%), and 'sadness' (25%). Notably, 'joy' emerged as the most dominant emotion expressed across the discussions, followed by 'anger' and 'sadness.'

RESULTS

➤ Sentiment Analysis Insights:

The sentiment analysis was conducted across various subreddits, notably "BlackLivesMatter," "stopasianhate," and "immigration," employing multiple techniques like VADER lexicon and BERT models. The analysis delved into the polarity of sentiments expressed in comments and posts within these forums.

The distribution of sentiments across the studied subreddits exhibited a diverse landscape of emotions. The findings revealed a multifaceted spectrum of sentiment expressions, comprising positive, negative, and neutral emotions. This analysis also highlighted evolving sentiment patterns over time, shedding light on the dynamic nature of discussions within these social movements.

The sentiment analysis offered valuable insights into prevalent themes, emotions, and sentiment shifts within these online discussions. It provided a comprehensive understanding of the emotional fabric underlying the discourse, portraying a rich tapestry of societal sentiments.

➤ Emotion Classification Proficiency:

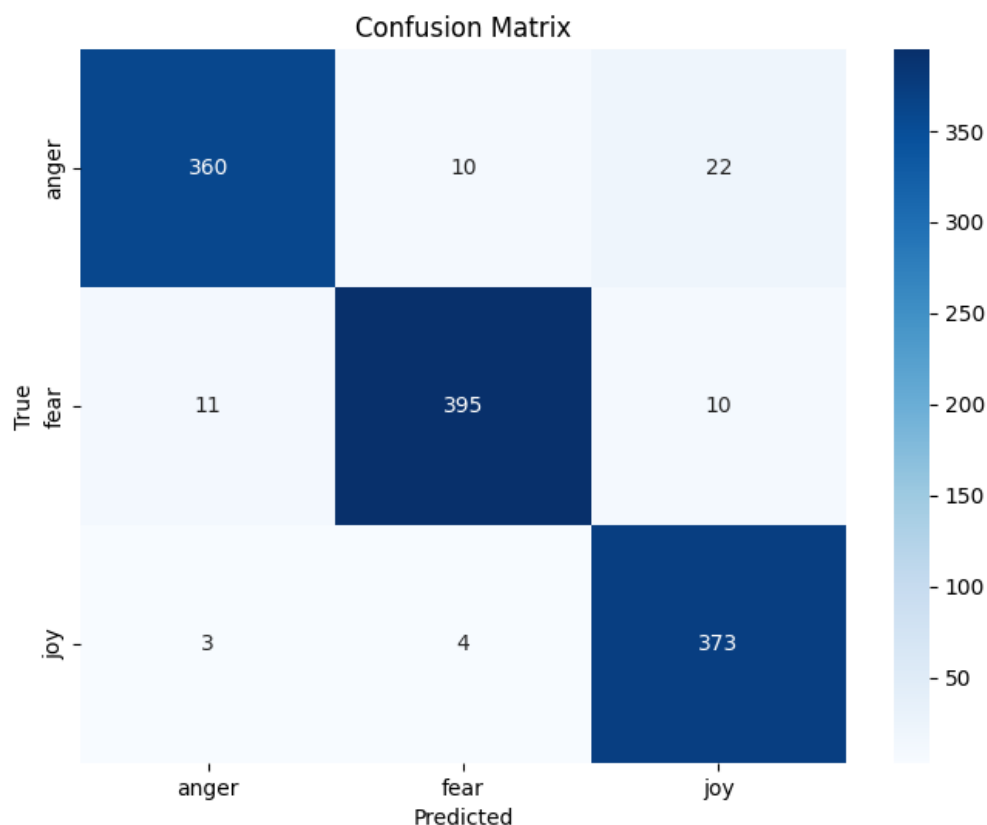
The LSTM-based emotion classification model, leveraging GloVe embeddings, showcased impressive proficiency in deciphering nuanced emotional expressions embedded within textual data. This model effectively categorized a range of emotions present in social media conversations related to socio-political movements.

The accuracy and effectiveness of the LSTM model underscored its ability to discern subtle emotional nuances, highlighting the complex emotional landscape inherent in online discussions. The model's proficiency in identifying and categorizing emotions further reinforced the significance of AI-driven approaches in understanding and analyzing societal emotions.

➤ Confusion Matrix:

This confusion matrix helps in understanding how well the model performs in identifying specific emotions such as "Joy," "Fear," and "Anger." By analyzing

these values, you can gauge which emotions the model might misclassify more often and which it predicts accurately, offering insights into the strengths and weaknesses of the emotion classification model within the context of your project.



- The row represents the actual / true emotions.
- The column signifies the predicted emotions.
- Values in the diagonal (from top left to bottom right) represent correct predictions (True Positives).
- Values off the diagonal represent incorrect predictions (False Positives and False Negatives).

This matrix allows for a clear visualization of the model's performance in predicting different emotions. The higher the values on the diagonal, the better the model's accuracy in predicting those specific emotions.

CITATIONS

1. [https://ieeexplore.ieee.org/abstract/document/8100646?casa_token= RM_FedpHmRgAAAAA:P5Dyo_wKpCNpBljvsq922yq1QziibzOmRZOXQb_YWgifqNEaB8j0UAACGTvuu6v9hTVcPJlivg](https://ieeexplore.ieee.org/abstract/document/8100646?casa_token=_RM_FedpHmRgAAAAA:P5Dyo_wKpCNpBljvsq922yq1QziibzOmRZOXQb_YWgifqNEaB8j0UAACGTvuu6v9hTVcPJlivg)
2. <https://pubs.aip.org/aip/acp/article/2916/1/030009/2926280/Using-social-media-to-predict-social-unrest-A>
3. <https://www.degruyter.com/document/doi/10.1515/opis-2022-0141/html?lang=en>
4. <https://www.elibrary.imf.org/view/journals/001/2021/263/article-A001-en.xml>
5. <https://arxiv.org/ftp/arxiv/papers/1805/1805.00358.pdf>
6. <https://b2find.dkrz.de/dataset/ac748c64-dc48-5572-ae86-1705251bcd95>