

Benchmarking Machine Learning Techniques in Under-Resourced Contexts: Analysis of Public Perceptions of Government Policies from South Karnataka Reddit Discourse

Varnica Sharma

Dept. of Applied Statistics and
Data Science, Prasanna School of
Public Health Manipal Academy of
Higher Education

Manipal, Karnataka, India
varnicasharma1515@gmail.com

Aman Tripathi

Dept. of Applied Statistics and
Data Science, Prasanna School of
Public Health Manipal Academy of
Higher Education

Manipal, Karnataka, India
amantripathi3273@gmail.com

K.M. Kavitha

Dept. of Applied Statistics and
Data Science, Prasanna School of
Public Health Manipal Academy of
Higher Education

Manipal, Karnataka, India
karimbikavitha@gmail.com

Abstract—Understanding public perception is crucial for the successful implementation of government policies. This study employs established Natural Language Processing (NLP) techniques to gauge public sentiment of government policies in Indian context and identify key themes related to policies in the South Karnataka region of India. Traditional machine learning techniques such as Support Vector Machine (SVM), Decision Tree (DT), Naïve Bayes (NB), K-Nearest Neighbors (KNN) and Random Forest (RF) along with deep learning techniques such as DistilBERT were explored in uncovering public opinion on government policies under very low training data conditions. Latent themes in the public discourse were identified using K-Means Clustering and BERTopic. Our findings reveal a clear performance trade-off for the highly nuanced task of perception mining, as all the chosen models under-performed under low-data settings with SVM emerging as the best performing model with a 41.5% accuracy slightly higher than the deep learning model DistilBERT with an accuracy of 40.4%. For thematic clustering, K-Means produced highly coherent and separable topics with 94.6% accuracy, demonstrating its effectiveness in identifying keyword-driven themes, whereas BERTopic identified more semantically nuanced topics that were less distinct.

Keywords—public perception mining, government policy analysis, low-resource NLP, reddit discourse analysis, sentiment analysis, machine learning benchmarking, DistilBERT, K-Means clustering, topic modeling

I. INTRODUCTION

Effective governance in the public sector hinges on the alignment between policy design and public acceptance. In the digital age, social media platforms have emerged as invaluable, albeit noisy, sources of unfiltered public opinion. This study harnesses the semi-anonymous, discussion-oriented platform Reddit to analyze citizen discourse surrounding Indian government policies in South Karnataka, India, a region characterized by rapid urbanization and evolving public service demands. We conduct a comparative study of two distinct Machine Learning (ML) paradigms: the traditional one which relies on statistical text features, and deep learning leveraging contextual embeddings for uncovering public perceptions on governmental policies. By evaluating the versatility of chosen models on the dual tasks of emotion classification and thematic clustering, we aim to provide a nuanced understanding of public sentiment and offer

a critical perspective on the models employed for such analysis exclusively in under-resourced contexts of public policy perception under very low training data conditions.

The study presented in this paper might be viewed as a benchmarking application of established machine learning and deep learning techniques to a novel, under-resourced data context: public discourse on government policies originating from Reddit users in South Karnataka, India. Our work does not introduce novel machine learning models, but rather provides a crucial proof-of-concept for the effective deployment of current methodologies in analyzing a politically and geographically distinct digital community. Existing literature on political opinion mining in India has predominantly focused on large-scale national sentiment or metropolitan areas. There is significant gap in studies that isolate and analyze the opinions of regional, linguistically diverse communities such as those in South Karnataka. Reddit discourse provides a level of anonymity and depth often absent in other social media platforms, potentially offering unfiltered views on sensitive government policies. The primary research gap this study addresses is the absence of an established machine learning performance baseline for analyzing public discourse in this under-resourced, regional context of public policy perception. While global benchmarks exist, their utility in predicting sentiment within a unique, low-resource dataset of public policy perception remains untested. The novelty of this research lies in its empirical contribution—establishing the first comparative performance benchmark of machine learning models for topic modeling and sentiment analysis of public policy perception, a localized dataset of South Karnataka Reddit comments.

II. MOTIVATION

A government policy's success is determined by the target citizens' perception, reception, and realized benefit, apart from the quality of its design and implementation. Indian public policy is characterized by a crucial disparity between official policy intent and the prevailing public discourse concerning its operational success and real-world impact. This study is an attempt to identify and analyze the divergence between the formal goals and desired outcomes of a policy (policy intent) and the way it is discussed, interpreted, and understood by the public (public discourse). In the said context, the contributions of the work presented in the paper include:

- Showcase how a combination of established ML and NLP techniques (SVM, Term Frequency-Inverse Document Frequency (TF-IDF), etc.) can be used to derive meaningful insights from a unique social media dataset on Indian policies.
- Provide a comparative analysis of how existing techniques (SVM vs. DistilBERT) can be effectively applied to a specific, underexplored dataset on Indian policies, thus serving as a useful benchmark for future research in this area.
- Gain new knowledge and understanding about public discourse on Indian policies, which was previously unavailable.
- Furthermore, these findings offer a critical, low-resource performance baseline for policymakers and NLP practitioners working on regional Indian public opinion.

III. REVIEW OF NLP FOR POLICY AND SENTIMENT ANALYSIS

Analysis of public sentiments across various domains has evolved significantly over the past two decades. Early sentiment analysis methods used handcrafted features to represent text. Models like SVM and NB relied on techniques such as n-grams and TF-IDF, but struggled to capture complex linguistic nuances like sarcasm and contextual meaning [1]. Specifically in the health domain, studies employing Twitter data to monitor public sentiment during the COVID 19 pandemic, often observed that public emotion shifted rapidly in response to government announcements and health outcomes [1].

The rise of deep learning, particularly transformer-based models like BERT, has revolutionized sentiment analysis by providing a richer understanding of semantic context [2]. Unlike earlier methods, these models could interpret complex arguments and subtle opinions that are not explicitly positive or negative [3]. Topic modeling techniques, such as the transformer-based BERTopic, have also emerged as powerful tools for identifying key themes in large corpora without supervision [4].

A major challenge in public policy is the disconnect between official goals and on-the-ground realities. For example, policies aiming for universal immunization may overlook regional infrastructure disparities, creating implementation barriers and resource mismatches. This gap is common in low and middle-income countries [World Health Organization (WHO) 2015], highlighting the need for methods that can capture grassroots sentiment. Similarly, analyses of education policies, such as India's National Education Policy (NEP) 2020, reveal gaps between goals and challenges (e.g., insufficient teacher training) [5]. For instance, an analysis of NEP 2020 Twitter data identified shifting public sentiment [6]. These parallels underscore the need for methodologies that help understand intent-implementation disparities across sectors [WHO 2015].

Existing literature demonstrates that NLP methodologies, validated in education and health policy contexts, provide a robust framework for bridging intent-implementation gaps. Synthesizing social media data offers policymakers a dynamic toolset for responsive, equity-focused reform [7]. Addressing cultural and demographic biases further enhance the model's

scalability across diverse public health landscapes [WHO, 2015] [8]. A comparative study using machine learning for policy-related sentiment reinforces the value of benchmarking different models on a specific, topical dataset [9]. A highly relevant study employing BERT-based model to analyze public sentiment on a specific policy issue in the Indian context highlights the value of using modern NLP techniques on social media data to inform policymakers [10].

While a significant body of NLP research exists for analyzing political discourse in English and other major languages, there's a growing need to apply these techniques to under-resourced contexts. This includes analyzing social media discourse in specific regions, such as the focus on South Karnataka in this study, and dealing with challenges like code-mixing (English and regional languages) and limited data availability. There is a shift towards hybrid methods that combine the strengths of different techniques. For example, unsupervised topic models like BERTopic are used to identify key themes, which can then be fed into supervised models for detailed sentiment or classification tasks. This allows for both broad exploratory analysis and specific, targeted insights. Our study builds on the prior works [1, 4, Table I] by providing a direct comparative analysis in the under-resourced context of regional Indian policy discourse.

IV. METHODOLOGY

To understand the public perception on governmental policies we initially began by procuring the data set sourced from Reddit. To uncover public opinion on available policies, classification was performed followed by clustering for identifying latent themes. The workflow illustrating the steps from data collection through classification and identifying themes is summarized in (Fig. 1).

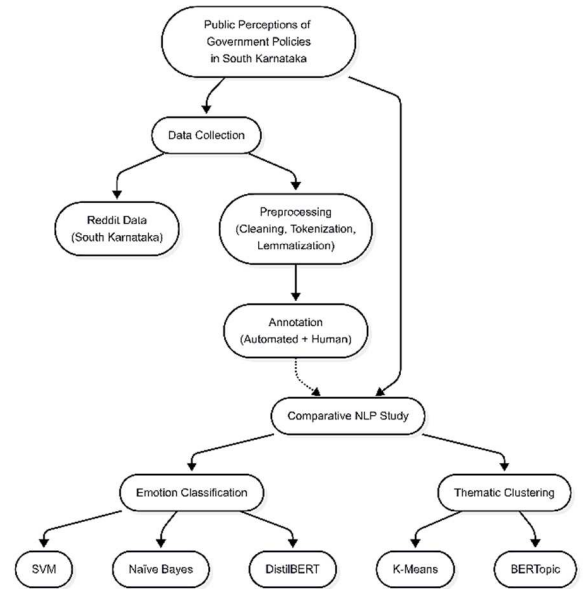


Fig. 1. Workflow illustrating data collection through analysis.

A. Text Acquisition and Corpus Building

To procure the dataset on public discourse on Indian policies, 735 Reddit comments were scraped using the PRAW API¹ from Reddit. The final CSV file consisted of five columns including “Comment Text”, “Article Title”, “Article URL”, “Article Publication Date”, and “Region”. The region-specific responses were derived from subreddit, with region specific scraping centered on South Karnataka². The

For the unsupervised cluster analysis, the afore-said dataset was merged with a corpus of 461 comments with search specifically focusing on public health policies to form an expanded dataset of 1196 comments, increasing its diversity and robustness. The comments scraped from selected YouTube videos and Reddit posts majorly focused on public health policies. Over 500 comments were initially scraped and then filtered to remove spam, emojis, links, and irrelevant text using natural language pre-processing techniques. Analyzing the dataset we observe that the word “government” is the most frequently mentioned, indicating that users place primary responsibility for healthcare outcomes on the state.

The graph displays the frequency of the top 20 most common words across 20 different samples. The word 'government' is the most frequent, appearing in approximately 95 samples. The frequency drops significantly for 'people' (approx. 76) and 'hospital' (approx. 74). The word 'scheme' is the 10th most frequent, appearing in about 71 samples. The frequency continues to decline, with 'year' at 57, 'govt' at 56, 'food' at 54, 'health' at 48, 'india' at 47, 'good' at 45, 'state' at 45, 'money' at 44, 'card' at 43, 'tax' at 37, 'private' at 33, 'pay' at 33, 'salary' at 32, 'ayushman' at 32, 'insurance' at 31, and 'citizen' at 31.

Sample	Word	Count (approx.)
1	government	95
2	people	76
3	hospital	74
4	scheme	71
5	year	57
6	govt	56
7	food	54
8	health	48
9	india	47
10	good	45
11	state	45
12	money	44
13	card	43
14	tax	37
15	private	33
16	pay	33
17	salary	32
18	ayushman	32
19	insurance	31
20	citizen	31

B. Text Pre-processing

C. Segregation of Public Perceptions

1) *Lexicon-based Emotion Labeling*: We first generated emotion labels for the 735 comments in the corpus using the lexicon-based NRClex library³. This unsupervised step leverages a predefined dictionary of emotional terms. This enabled each comment a dominant emotion label “anger”, “joy”, “sadness”, “fear”, “surprise” and so forth. 80% of the labeled dataset was used as training set and 20% as test set.

3) *Classification*: We trained and evaluated five traditional machine learning models such as SVM, DT, NB, KNN and RF. In parallel, we evaluated a pre-trained transformer model DistilBERT⁴ [12] via HuggingFace’s pipeline to label each comment’s emotion.

Recognizing the limitations of pre-trained models on specialized data, we established a rigorous validation protocol. A human-annotated ground truth dataset was created by having three human raters independently label a random sample of 100 comments. Based on inter-annotator agreement “true” emotion for each comment was determined by a majority vote. The performance of the Large Language Model (LLM) baseline (DistilBERT) was then quantitatively

validated against this 100-comment human gold set (three-rater majority labels), followed by a qualitative error analysis to understand its specific strengths and weaknesses.

In addition to the segregation of comments based on perceptions, we further considered identifying the latent themes underlying the perceptions to understand stakeholder opinions on governmental policies. On the merged dataset of 1196 comments, we compared two clustering techniques to identify latent topics in the discourse. We explored K-Means and BERTopic [13] models and are as detailed in the subsections that follow.

2) *BERTopic Modeling*: As a modern alternative, we used BERTopic5, a transformer-based model that clusters comments based on the semantic similarity of their sentence embeddings, allowing for the discovery of more nuanced topics. Using BERTopic we extracted the top keywords of each topic to characterize main themes.

A. Emotion Classification Results

TABLE I. PERFORMANCE OF CLASSIFICATION MODELS.

Following prediction, the LLM’s resulting distribution is shown in (Fig. 5) , contrasting sharply with the reference ground truth. The DistilBERT model, when evaluated against the human-annotated ground truth, achieved a validated accuracy of 40.4%. The pre-trained LLM did not significantly outperform the simpler traditional method suggesting that without domain-specific fine-tuning, even powerful general models struggle to interpret regional discourse.

Fig. 4. Emotion Distribution from NRCLex Labels (Original Dataset).

Fig. 5. Emotion Distribution from LLM Predictions (Original Dataset).

A qualitative error analysis, detailed in Table 1, provides insight into the LLM’s behavior. The model correctly identified straightforward expressions of “confusion” and “gratitude”. However, it faltered on comments where negative sentiment was conveyed indirectly, misinterpreting a comment expressing “disappointment” about water shortages as “confusion”, and a comment expressing “anger” at government failure as “disappointment”.

Evaluation (Comment Text)	Actual Emotion	Predicted Emotion
I'm not sure what to think anymore...	confusion	confusion
I'm so grateful for the people who are...	gratitude	gratitude
Karnataka has highest rainfall udupi yet...	disappointment	confusion
This government is gone bonkers...	anger	disappointment

B. Thematic Clustering Results

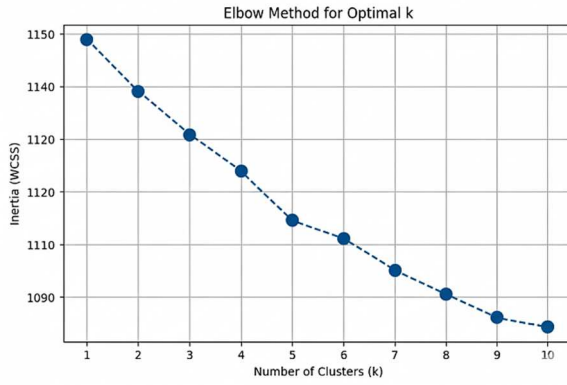


Fig. 6. Elbow plot justifying $k=5$ for K-Means.

Fig. 7 shows the distribution of documents across each of the 5 clusters identified using the K-Means clustering.

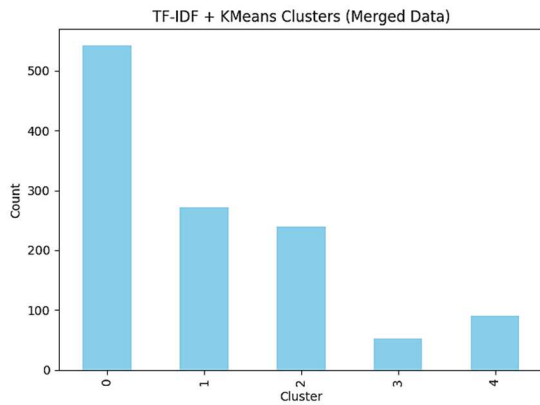


Fig. 7. Distribution of documents across the 5 K-Means clusters.

Table III summarizes the top 25 unigrams and selected bigrams identified for each of the 5 clusters using the K-Means Clustering. Table IV depicts the top terms identified using the BERTopic technique.

TABLE III. TOP TERMS FROM K-MEANS CLUSTERING.

Cluster	Top Terms
0	Unigrams: health, good, food, people, govt, scheme, like, india, hospitals, healthy, card, deleted, getting, money, great, time, don, protein, know, mandya, need, yes, scss, ayushman, salary. Bigrams: -
1	Unigrams: government, people, just, state, tax, central, like, money, don, india, work, companies, way, pay, private, need, make, good, point, bjp, years, city, karnataka. Bigrams: central government, state government
2	Unigrams: karnataka, kannada, language, state, people, just, north, like, hindi, south, bangalore, states, english, want, mysore, don, india, yes, learn, culture, bjp, land, hate, govt. Bigrams: north karnataka
3	Unigrams: mysuru, mysore, airport, city, don, posts, bengaluru, rich, big, years, real, post, new, young, bangalore, people, great, agree, golden, details, decade, like, develop, money, development. Bigrams: -
4	Unigrams: bengaluru, city, karnataka, urban, bangalore, https, namma, post, traffic, south, try, www, com, road, come, days, just, kannada, india, rural, did, places. Bigrams: namma bengaluru, karnataka bengaluru, https www

TABLE IV. TOP TERMS FROM BERTOPIC.

Cluster	Top Terms
0	karnataka, state, bjp, people, guy, media, tn, government, just, north, states, bihar, districts, trying, away, politicians, didnt, like, ka, news, good, auto, bangalore, development, india
1	hospitals, hospital, private, doctors, patient, card, ayushman, people, treatment, insurance, healthcare, patients, government, money, private hospitals, medical, health, blood, dont, ayushman card, know, covered, scheme, pay, india
2	food, protein, healthy, eat, indian, foods, sugar, diet, products, health, good, cooked, nutrition, india, eggs, indians, high, meat, stomach, fruits, rice, like, daily, cook, carbs
3	language, kannada, hindi, learn, karnataka, english, languages, south, speak, spoken, like, dont, kodava, indian, just, culture, south indian, people, tongue, tulu, professional, public, mother tongue, learning, mother
4	government, protests, just, protest, doing, lapse, fuck, work, action, question, instead, security lapse, fuck government, public, central, president, local government, pockets, people, politicians, central government, protesting, elected, power

While testing the coherence of the K-Means clusters by training an SVM to predict them, the model achieved an exceptionally high accuracy of 94.6%. This suggests that the clusters are linearly separable and thematically consistent. The PCA visualization in (Fig. 8) further illustrates these distinct groupings. In contrast, the SVM accuracy on BERTopic labels was much lower at 52.5%.

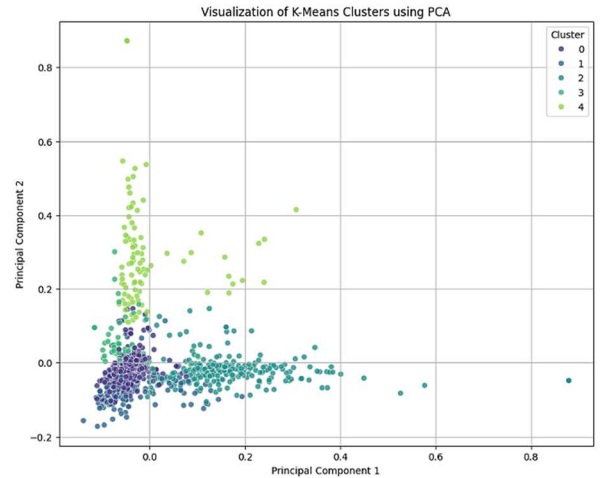


Fig. 8. Visualization of K-Means clusters using PCA.

C. Limitations and Constraints

Our work solely relies on analysis of public perceptions sourced from online discourses. This might underrepresent demographics with limited internet access or digital literacy, such as older adults, low-income individuals, and those in rural areas. Their perception of public policies, particularly those affecting healthcare, are crucial, nevertheless might go unheard. As data acquisition was sourced from South Karnataka and specific subreddits, it might not accurately reflect the opinions of the entire region or country. As the scraped comments were predominantly in English, with occasional local language terms, the analysis might not reflect the perceptions of majority of population communicating in regional languages. The use of Reddit and YouTube as data sources limits our analysis to comments from individuals who

actively use these platforms, which may not be a representative cross-section of society. The high frequency of terms, such as, “government” and “hospital” might reflect the concerns of an already engaged audience rather than the general population's primary issues. Thus, the dataset possibly overrepresents extreme or highly engaged viewpoints, while a silent majority with moderate or indifferent opinions is ignored, indicative of self-selection bias.

VI. CONCLUSION

This study conducted a comparative analysis of Machine Learning methods to dissect public discourse on health policies in South Karnataka, leveraging a unique dataset of approximately 1,200 comments from South Karnataka Reddit discourse. In an under-resourced, very low training data conditions, primary success of our work is measured by two key contributions: the rich thematic insights uncovered and the performance benchmark established for machine learning application in this regional setting.

Our validated results yielded a key insight: a powerful, pre-trained DistilBERT model (40.4% accuracy) performed comparably to a much simpler lexicon-based SVM (41.5% accuracy) for emotion classification, underscoring the critical importance of domain-specific validation. For thematic analysis, we found that K-Means produced highly coherent and separable clusters (94.6% SVM accuracy), indicating that keyword-based topics were more distinct than the semantically richer but overlapping topics from BERTopic.

By applying and comparing established machine learning and deep learning techniques such as SVM, DT, NB, KNN, RF, and DistilBERT we have established the first performance baseline for sentiment analysis and topic classification of public discourse on health policies within the South Karnataka Reddit corpus. While we acknowledge that our study is constrained by the small dataset size, this benchmark is essential for any subsequent, larger-scale studies in this domain, providing a necessary proof-of-concept for the feasibility of these methods in a low-resource context. While a pre-trained model like DistilBERT is effective, its performance can be significantly enhanced by fine-tuning it on the specific dataset. Fine-tuning a transformer model like BERT on this local, domain-specific data is a critical next step that could significantly improve emotion and topic classification accuracy, considering the feasibility of larger corpus of Indian public policy discourse.

While local language terms were occasionally observed in our corpus, with multilingual embeddings, employing multilingual models like SLM-RoBERTa, IndicBERT or Mbert [14], pre-trained on text from multiple languages including a wide range of Indian languages, would allow the model to more accurately process comments that mix English and Kannada or other regional languages rather than restricting to English parts. Cross-regional comparison would

further enable us to generalize our findings to vast majority of the population in India, which primarily communicates in regional languages.

We envision a more robust system that is both interpretable and highly accurate by adopting a hybrid framework that merges the structured, keyword-based insights from traditional methods like TF-IDF with the deep contextual understanding of transformer models. Extending our analysis with multimodal data by including images or news links referenced by comments would further improve reliability and depth of public opinion insights, helping tailor communication strategies and policies more effectively. Temporal analysis capturing how sentiment and topics evolve over time with the new policy announcements is yet another potential area for exploration.

REFERENCES

- [1] K. L. Tan, C. P. Lee, and K. M. Lim, "A survey of sentiment analysis: Approaches, datasets, and future research," *Applied Sciences*, vol. 13, no. 7, p. 4550, 2023.
- [2] L. Zhang, S. Wang, and B. Liu, "Deep learning for sentiment analysis: A survey," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 8, no. 4, p. e1253, 2018.
- [3] G. Aadil and D. Samad, "Traditional or deep learning for sentiment analysis: A review," *Multidiszciplináris tudományok*, vol. 12, no. 1, pp. 3–12, 2022.
- [4] M. Esh, "Sentiment analysis in chatgpt interactions: Unraveling emotional dynamics, model evaluation, and user engagement insights," *Technical Services Quarterly*, vol. 41, no. 2, pp. 160–174, 2024.
- [5] K. C. Karthikeyan C, "Emerging challenges in implementation and practice of nep (national education policy) for teachers, students and stakeholders," Jan. 2021.
- [6] J. Salam, "Draft national education policy (nep), 2019 and jingoistic nationalism," *The People's Chronicle*, vol. VII, 2019.
- [7] V. Anoop and S. Sreelakshmi, "Public discourse and sentiment during mpox outbreak: an analysis using natural language processing," *Public Health*, vol. 218, pp. 114–120, 2023.
- [8] X. Gong and J. Wen, "The text analysis of national education policy 2020," *Journal of Education, Teaching and Social Studies*, vol. 5, p. p95, Feb. 2023.
- [9] D. S. Zainulabdeen, M. Çevik, and M. M. Abdulrazzaq, "A comparative study of classification algorithms for sentiment analysis of covid 19 vaccine opinions using machine learning," in *2024 International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)*. IEEE, 2024, pp. 1–6.
- [10] P. Khot, M. Vishwakarma, and V. Shukla, "Bert-based sentiment analysis of indian covid-19 tweets for policy making," *Iconic Research And Engineering Journals*, vol. 7, no. 8, pp. 26–33, 2024.
- [11] O. Ütük Bayılmış and S. Orhan, "Decoding digital labor: A topic modeling analysis of platform work experiences," *MDPI*, vol. 13, no. 9, p. 819, 2025.
- [12] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *ArXiv*, vol. abs/1910.01108, 2019.
- [13] M. Grootendorst, "Bertopic: Neural topic modeling with a class-based tf-idf procedure," *arXiv preprint arXiv:2203.05794*, 2022.
- [14] R. Pathak and A. Kumar, "Aspect-based sentiment analysis in hindi language by ensembling pre-trained mbert models," *Electronics*, vol. 10, no. 21, p. 2641, 2021.