

Ahongsangbam Dorendro^{1*},H. Mamata Devi²

Challenges in Determining the Authenticity, Honesty, and Intentions of Opinions Expressed on Twitter and Sentiment Analysis



Abstract. Social media platforms, notably Twitter, have emerged as pivotal arenas for the dissemination of opinions spanning diverse domains, with politics standing prominently among them. However, discerning the authenticity, honesty, and underlying intentions behind these expressed opinions poses a formidable challenge. This study delves into the multifaceted obstacles associated with gauging the credibility of opinions proliferating on Twitter, with a particular focus on the pervasive influence of social dynamics, the manipulation of narratives through the creation of new accounts, and the heightened complexities surrounding political discourse. Social dynamics wield considerable influence over the credibility of opinions expressed on Twitter. Users often engage in echo chambers, where confirmation bias amplifies certain viewpoints while stifling dissenting voices. Additionally, the phenomenon of social influence leads individuals to conform to prevailing opinions, further complicating the task of discerning genuine sentiment. Moreover, the proliferation of fake accounts designed to manipulate narratives poses a significant challenge to opinion credibility assessment. To address these challenges, novel approaches must be incorporated into Twitter data Sentiment Analysis frameworks.

Keywords: Twitter, Social Dynamics, Echo Chambers, Astroturfing, Sentiment Analysis.

1 Introduction

1.1 Background and related works

In recent years, social media platforms, spearheaded by Twitter, have undeniably transformed the landscape of public discourse. They offer individuals an unprecedented avenue to share their thoughts, ideas, and perspectives with a global audience instantaneously. The allure of these platforms lies in their ability to democratize expression, breaking down geographical barriers and allowing voices from all corners of the world to be heard. However, beneath the surface of this digital revolution lies a complex web of challenges that threaten the authenticity and integrity of the opinions shared.

One of the most glaring issues is the obscured nature of authenticity, honesty, and intentions behind the opinions proliferating on social media. While users may express themselves freely, the motivations driving their words can often remain veiled. The anonymity afforded by the internet, coupled with the ease of creating pseudonymous accounts, further complicates matters, blurring the lines between genuine expression and manipulation.

In this environment, deciphering the true meaning and intent behind social media posts becomes increasingly challenging. It is crucial to acknowledge the limitations of sentiment analysis, particularly its inability to fully capture the nuances of human communication. Without addressing the complexities inherent in assessing authenticity, honesty, and intentions on social media, conclusions drawn solely from sentiment analysis risk painting an incomplete picture. While sentiment analysis can provide valuable insights into the prevailing mood of online conversations, it cannot definitively ascertain the underlying motives driving those sentiments [1].

Thus, as we delve deeper into the realm of social media discourse, it becomes imperative to approach our analysis with a critical eye, recognizing the inherent challenges and limitations at play. By doing so, we can strive to foster more meaningful and informed discussions in an increasingly digital world.

Sentiment analysis, also referred to as opinion mining, is a computational technique aimed at identifying the sentiment expressed within a piece of text, such as social media posts, reviews, or feedback [2]. Its primary objective is to gauge the overall attitude or opinion of the author or the public towards a particular subject, be it a person, object, event, or topic. The proliferation of digital content on the internet has led to an exponential growth in data availability, making sentiment analysis an indispensable tool for extracting meaningful insights from this vast reservoir of information. By employing natural language processing (NLP) and machine learning algorithms, sentiment analysis algorithms can automatically categorize text into positive, negative, or neutral sentiments, or even provide a graded score indicating the strength of sentiment. The application of sentiment analysis is diverse and spans various domains. In the realm of business and marketing, it serves as a crucial tool for analyzing customer feedback, reviews, and social media conversations surrounding products and services. By understanding the sentiments expressed by consumers, companies can adapt their strategies accordingly, whether it involves refining their products, addressing customer concerns, or leveraging positive feedback for marketing purposes.

Moreover, sentiment analysis plays a pivotal role in political discourse, especially in the age of social media [3]. Platforms like Twitter have become prominent arenas for political discussions and debates, with politicians, parties, and citizens alike expressing their opinions and viewpoints. Analyzing sentiment in social media posts can provide valuable insights into public sentiment towards political parties, candidates, policies, and current events. This information is instrumental for political campaigns, policymaking, and understanding public opinion dynamics

¹Department of Computer Science, Manipur University, Imphal West, India. Email: ah.dorendro@gmail.com, [0000-0003-3941-8685]

²Department of Computer Science, Manipur University, Imphal West, India. Email: mamata_dh@rediffmail.com

1.2 Objective

This study aims to identify the key challenges in determining the authenticity of opinions expressed on Twitter. Authenticity, in this context, refers to the genuineness or sincerity of the expressed sentiments. It's a critical aspect because understanding authenticity is essential for making accurate interpretations of social media data, especially when drawing conclusions or making decisions based on that data. The study will further elaborate on the impact of social dynamics on shaping and influencing opinions. Finally, we will identify potential strategies to minimize the impact of inauthentic and insincere content in our sentiment analysis dataset. These strategies will be integrated into the text preprocessing and feature selection stages during the development of our Sentiment Analysis model framework.

2 Key challenges in determining the authenticity of opinions expressed on Twitter

Determining the authenticity of opinions expressed on social media platforms is a complex challenge influenced by various social dynamics and opinion formation processes. The challenges are compounded by the creation of fake accounts and manipulation tactics. Some of the key challenges are discussed below.

2.1 The Echo Chamber Effect

One significant challenge in assessing the authenticity of opinions on Twitter is the echo chamber effect [4]. Users tend to follow accounts and engage with content that aligns with their existing beliefs, creating an insular environment where dissenting opinions are often ignored or dismissed. This echo chamber effect can lead to the reinforcement of existing biases and the distortion of public discourse [5].

Consider a hypothetical scenario where there's a contentious political issue, such as climate change. Users on Twitter often follow accounts and engage with content that reflects their pre-existing beliefs. Those who strongly believe in the urgency of addressing climate change are likely to follow environmental organizations, scientists, and activists who share articles, data, and opinions supporting their stance.

On the other hand, users skeptical about climate change may follow accounts of think tanks, politicians, or individuals who share content questioning the scientific consensus on climate change. As a result, each group is exposed primarily to content that reaffirms their beliefs, creating an echo chamber effect.

Within these echo chambers, users may share, retweet, and amplify information that aligns with their views while dismissing or ignoring dissenting opinions. This selective exposure can lead to the proliferation of misinformation, as users are less likely to critically evaluate information that challenges their beliefs.

For instance, in the case of climate change, those who deny its existence might share articles from fringe websites or pseudoscientific sources that cast doubt on established climate science. Conversely, those who advocate for action on climate change might amplify reports from reputable scientific institutions and news outlets that support the consensus view.

As a result, the echo chamber effect amplifies polarization and stifles constructive dialogue. It becomes increasingly challenging for users to engage in meaningful discussions across ideological divides, hindering the potential for consensus-building and informed decision-making.

In this way, the echo chamber effect on Twitter can have profound implications for public discourse, perpetuating misinformation, reinforcing biases, and impeding the exchange of diverse perspectives.

2.2 Virality and Herd Mentality

Virality and Amplification. The viral nature of content on Twitter contributes to the bandwagon effect, where opinions gain traction simply due to their popularity rather than their merit [6]. Twitter's design, with its emphasis on brevity and rapid dissemination of information, lends itself well to the viral spread of content. Tweets that are catchy, controversial, or emotionally charged are more likely to be shared widely, often reaching thousands or even millions of users within a short span of time. This rapid amplification can turn relatively obscure opinions or ideas into widespread trends. When content goes viral, accumulating a large number of likes, shares, and comments, it provides a sense of *social validation*. Users perceive popular content as more credible, trustworthy, and worthy of attention simply because many others have engaged with it. This social validation creates a feedback loop where the popularity of content reinforces its perceived value, leading more users to join in and contribute to its virality.

Herd Mentality. The herd mentality describes the tendency for individuals to conform to the behavior or opinions of a larger group, often without critically evaluating those opinions themselves [7]. On platforms like Twitter, where metrics like number of likes, retweets, and comments serve as indicators of popularity and social approval, users may be inclined to follow the crowd rather than formulating their own independent judgments.

In the context of social media, viral content creates a virtual bandwagon that attracts users to jump on board and participate in discussions or activities related to that content. Users may be swayed by the sheer number of likes, retweets, or comments, leading to the adoption of opinions without critical evaluation. Users often use these metrics as cues for deciding which content to engage with or endorse. The challenge posed by the herd mentality on social media is that opinions and behaviors may be driven more by social pressure and the desire for social approval than genuine belief or critical evaluation. Users may hesitate to express dissenting opinions or engage in nuanced discussions if they perceive that the majority holds a different viewpoint. This obscures the authenticity of opinions expressed over social media platforms like twitter, as they may be driven more by social pressure than genuine belief.

2.3 Astroturfing and Inauthentic Amplification

Astroturfing is a deceptive tactic where fake accounts or manipulated existing accounts are used to create the illusion of widespread grassroots support for a specific opinion, idea, or cause. Fake accounts on social media platforms refer to profiles that are created with the intent to deceive or manipulate others. Fake accounts are typically created for various purposes, including spreading misinformation, amplifying certain viewpoints, manipulating public opinion. This practice can significantly distort online discourse by blurring the line between genuine and orchestrated opinions. It becomes challenging for individuals to discern whether the support they see is authentic or artificially generated.

Inauthentic amplification often involves the use of bots and automated systems, which can further complicate matters. Bots are software programs designed to perform tasks automatically, such as posting messages, liking content, or following accounts. When used in astroturfing, bots can artificially amplify certain viewpoints by creating a facade of popularity or consensus [8]. They can make opinions that may not have significant genuine backing appear to be widely accepted or endorsed [9]. This manipulation of public perception can have far-reaching consequences, affecting everything from political debates to consumer choices on social media platforms. It underscores the importance of verifying information from multiple sources, especially in the digital age where misinformation and manipulation are prevalent.

2.4 Coordinated Campaigns and Manipulation Tactics

In addition to astroturfing, coordinated campaigns by organized groups or individuals with specific agendas pose a significant challenge [10]. Coordinated campaigns refer to organized and strategic efforts by individuals, groups, or entities to promote a particular agenda, narrative, or message on social media platforms. These campaigns often involve collaboration among multiple actors who work together to amplify content, target specific audiences, and manipulate online discourse. Coordinated campaigns can be carried out by various parties, including political organizations. These actors may employ sophisticated manipulation tactics, such as coordinated hashtag campaigns. Coordinated campaigns may also employ psychological tactics to manipulate emotions, biases, and cognitive vulnerabilities. This can include appealing to fear, anger, or tribal instincts, exploiting cognitive biases such as confirmation bias or the illusion of truth effect, or using persuasive techniques derived from behavioral psychology. By exploiting these psychological mechanisms, manipulators can subtly influence user opinion without individuals being consciously aware of it.

2.5 Organic Targeting and Polarization

Organic targeting. It refers to tactics such as utilizing hashtags, mentions, and the reply feature to engage with specific audiences and accounts beyond the follower list on the social media platform [11]. It is a natural way where content is disseminated and reaches specific audiences without paid promotion or advertising. Organic targeting can reinforce echo chambers by exposing users primarily to content that aligns with their existing beliefs and preferences. Users within echo chambers may perceive opinions that resonate with their beliefs as more authentic, overlooking diverse perspectives and critical evaluation.

Polarization. Polarization on the other hand, is the tendency for individuals to gravitate towards extreme positions on various issues. Social media platforms, with their algorithm-driven content delivery systems, inadvertently exacerbate polarization by amplifying content that is more likely to evoke strong emotional responses. This often leads to the formation of ideological bubbles where users are isolated from dissenting opinions, further reinforcing their own beliefs and perceptions [12]. Eventually, it can lead to a skewed perception of authenticity, as users may primarily encounter content that reinforces their existing opinions and biases. At the same time, social media users may inadvertently overlook or dismiss opinions that challenge their worldview, contributing to a lack of diversity and authenticity in online discourse

The combination of organic targeting and polarization [13] has significant implications for retrieving actual insights from social media posts. It can contribute to the spread of misinformation and which in turn mask the true intention of opinions and sentiments expressed on social media platforms.

3 Impact of social dynamics and Sentiment Analysis Model framework

We have seen above that a lot of social dynamics and manipulation factors are in play that affects the opinions that are shared on various social medias. We have seen how astroturfing make it difficult to differentiate between genuine and orchestrated opinions. There is a strong affinity between a social media users' exposure to misinformation and the degree to which they are perceived as ideologically conservative [13][14]. Coordinated campaigns and manipulation tactics, herd mentality contributes to the spread of misinformation over social media platforms [7][15].

Overcoming the challenges associated with determining the authenticity, honesty, and intentions behind opinions expressed on Twitter and incorporating measures to tackle them during sentiment analysis model framework involves considering multiple parameters and adopting a multifaceted approach. Given below are some key parameters or features and strategies to address these challenges:

3.1 Account Creation Date:

We can identify the account creation date as a feature for sentiment analysis model framework. Examining the account creation date is crucial to identify potentially inauthentic accounts or those created with the intention of manipulating

narratives that is fake accounts. Weightage can be assigned to the overall tweet sentiments based on the age of the account. A threshold can also be derived which in turn decides if the entire tweet sentiment is to be evaluated or discarded.

3.2 Engagement Metrics

Metrics viz. Likes, Retweets, Comments, posting frequency can provide a valuable insight on the user behavior and tweets. High engagement metrics can indicate the popularity of an opinion but may also indicate manipulation. Understanding user behavior based on the above metrics can help identify potential manipulation tactics, such as coordinated campaigns or artificial amplification, allowing sentiment analysis models to flag suspicious activity and prioritize authentic opinions

3.3 User Verification

Verified accounts are less likely to be associated with inauthentic behavior. Verification adds a layer of authenticity to a user's identity. Moreover, a verified account are less likely to be a bot or a fake account solely created for the purpose of manipulation, coordinated campaign and dissemination of wrong information. A weightage may be assigned to the tweet of a verified user such that a tweet from a verified user has more weightage of being truthful than that of an unverified user.

3.4 Network Analysis

Examining the connections between users can reveal patterns of manipulation, such as coordinated campaigns or astroturfing. We can exploit the follower and following networks of twitter user to investigate patterns of manipulation. This will also reveal any instance of herd mentality of the user and which in turn can be used to evaluate the authenticity of the opinions expressed over social media platforms like twitter.

3.5 Real-Time Monitoring:

Real-time monitoring of social media data can be incorporated in the Sentiment Analysis model framework to detect sudden shifts in sentiment, anomalies in user behavior, or emerging trends. This will enable sentiment analysis models to respond promptly to changing dynamics, identify potential manipulation tactics, and maintain the authenticity of sentiment analysis insights

3.6 Sentiment Analysis Model Framework

By considering the above discussed parameters and implementing corresponding model framework, sentiment analysis can provide a more transparent, authentic, and trustworthy analysis of the political discourse expressed over social media. The model framework will assign different weights to samples in the training data based on weightage of the individual parameters and overall sentiment tweets.

Preprocess the data by removing or down-weighting terms associated with biased sentiments. For example, identify and adjust the weight of certain words or phrases that are disproportionately associated with a particular sentiment.

4 Conclusion

The challenges in determining the authenticity, honesty, and intentions behind opinions expressed on Twitter are multifaceted. Social dynamics, manipulation tactics, and the unique characteristics of political opinions on the platform contribute to a complex landscape. Incorporating these parameters in the sentiment analysis model framework will give a better and authentic analysis of sentiment expressed over social media platforms. Other text classification model like opinion spamming, fake reviews detection can be incorporated to enhance the performance of the Sentiment Analysis model.

References

- [1] Rita, P., António, N. & Afonso, A. Social media discourse and voting decisions influence: sentiment analysis in tweets during an electoral period. *Soc. Netw. Anal. Min.* 13, 46 (2023). <https://doi.org/10.1007/s13278-023-01048-1>
- [2] Turney, and Peter: Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews. *Proceedings of the Association for Computational Linguistics*. Philadelphia, (2002). <https://doi.org/10.48550/arXiv.cs/0212032>
- [3] Tumasjan, A., Sprenger, T., Sandner, P., & Welpe, I. Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment. *Proceedings of the International AAAI Conference on Web and Social Media*, 4(1), 178-185. (2010). <https://doi.org/10.1609/icwsm.v4i1.14009>
- [4] Cota, W., Ferreira, S.C., Pastor-Satorras, R. et al. Quantifying echo chamber effects in information spreading over political communication networks. *EPJ Data Sci.* 8, 35 (2019). <https://doi.org/10.1140/epjds/s13688-019-0213-9>
- [5] Cinelli, M., De Francisci Morales G., Galeazzi, A., Quattrociocchi, W., Starnini, M. The echo chamber effect on social media. *Proc Natl Acad Sci U S A.* (2021). 118(9):e2023301118. doi:10.1073/pnas.2023301118
- [6] The Decision Lab Bandwagon Effect, <https://thedecisionlab.com/biases/bandwagon-effect>, last accessed 2024/04/21
- [7] Pavlović-Höck, D. N. Herd behaviour along the consumer buying decision process - experimental study in the mobile communications industry. *Digital Business*, 2(1), Article 100018. (2022). <https://doi.org/10.1016/j.digbus.2021.100018>

- [8] Zhang, Y.; Song, W.; Koura, Y.H.; Su, Y. Social Bots and Information Propagation in Social Networks: Simulating Cooperative and Competitive Interaction Dynamics. *Systems*.11, 210. (2023). <https://doi.org/10.3390/systems11040210>
- [9] Koc-Michalska, K., Klinger, U., Bennett, L., & Römmele, A. (Digital) Campaigning in Dissonant Public Spheres. *Political Communication*, 40(3), 255–262 (2023). <https://doi.org/10.1080/10584609.2023.2173872>
- [10] Kubin, E., & von Sikorski, C. The role of (social) media in political polarization: a systematic review. *Annals of the International Communication Association*, 45(3), 188–206 (2021). <https://doi.org/10.1080/23808985.2021.1976070>
- [11] Kountouri, F, and A Kollias: Polarizing publics in Twitter through organic targeting tactics of political incivility. *Front. Polit. Sci.* 5 (2023). <https://doi.org/10.3389/fpos.2023.1110953>
- [12] Mosleh, M., Rand, D.G. Measuring exposure to misinformation from political elites on Twitter. *Nat Commun* 13, 7144 (2022). <https://doi.org/10.1038/s41467-022-34769-6>
- [13] Gaultney, I. B., Sherron, T., & Boden, C. Political polarization, misinformation, and media literacy. *Journal of Media Literacy Education*, 14(1), 59-81 (2022). <https://doi.org/10.23860/JMLE-2022-14-1-5>
- [14] Su, M.-H., Suk, J., & Rojas, H. Social Media Expression, Political Extremity, and Reduced Network Interaction: An Imagined Audience Approach. *Social Media + Society*, 8(1). (2022). <https://doi.org/10.1177/20563051211069056>
- [15] Su, X., P. Li, and X. Zhu: The Influence of Herd Mentality on Rating Bias and Popularity Bias: A Bi-Process Debiasing Recommendation Model Based on Matrix Factorization. *Behav. Sci.* 13, no. 63 (2023). doi: 10.3390/bs13010063