# Evaluating random sampling bias in sentiment analysis of social media data

**Xuanting Xiong**

College of Arts, Sciences & Engineering, University of Rochester, Rochester, New York, United States, 14627

xxiong6@u.rochester.edu

**Abstract.** In this age marked by a wealth of information, the relevance of social networks has increased in a manner that is analogous to an exponential growth curve. Notably, the content that is shared on these platforms has the potential to act as a reflection of the emotional states that people are now experiencing. The importance of emotions is brought to light in the research presented here, which makes use of a technique based on a review of the relevant literature to analyze the problem of random sample bias and the effects that it has on sentiment analysis. It is possible to draw the conclusion, on the basis of the findings of the research, that the problem of random sample propensity is not a sporadic or insignificant one. In addition, the findings of the study indicate the presence of multiple types of prejudice. Because of the potential repercussions that could result from doing a distorted sentiment analysis, it is really necessary to keep your method focused.

**Keywords:** Social Media Data, Random Sampling Bias, Sentiment Analysis

## 1. Introduction

In the era of vast amounts of information, social networks have become integral to individuals' daily lives. The platform provides individuals with the opportunity to articulate their ideas, opinions, and emotions. Media posts not only serve as individual posts, but also provide significant identity-related information that might represent a collective group of individuals. Sentiment analysis has become a crucial tool for extracting valuable insights from the vast volume of unstructured data present on social media platforms. The potential uses of gauging popular opinion on corporate products and assessing citizen responses to new laws by federal governments are extensive. Social networking platforms provide a significant amount of immediate feedback on a wide range of topics. The task of accurately analyzing this enormous amount of information has inherent difficulties. One notable problem, which serves as the central emphasis of this research, pertains to the inherent bias created by random sampling in the examination of beliefs derived from social media data. The potential consequences of sampling bias can have significant ramifications, as they have the capacity to distort insights and lead to erroneous decision-making.

Understanding the intricacies of this inclination, its origins, and its potential impacts is crucial for any individual or entity relying on sentiment analysis in the contemporary digital environment. This study delves into the intricacies of random sample bias and its implications for sentiment analysis.

## 2. Literature review

Sentiment analysis, also referred to as opinion mining, encompasses the computer process of ascertaining the favorable, unfavorable, or neutral tone of a written article. Social network platforms play a significant role in facilitating the comprehension of the collective opinions, perspectives, and emotions of the general populace towards specific topics, entities, or trends [1]. The utilization of social media by approximately 3.6 billion individuals globally has significant implications for commercial strategies, policymaking, and sociocultural tendencies [2].

The analysis of beliefs presents inherent obstacles. The presence of information noise, which encompasses spam, irrelevant content, and web lingo, has the potential to impede the accuracy of belief detection. The significance of the circumstances in which declarations are made on social media is therefore of utmost importance. The connotation of a certain term might vary depending on the situation, leading to a favorable or unfavorable perception [3]. Moreover, the investigation of beliefs must confront the intrinsic subjectivity of human emotions and communications, so acknowledging that same statements may be interpreted differently by other individuals or algorithms.

Regarding the utilization of random sampling in belief analysis, numerous studies have been conducted to investigate its effects. The process of random sampling entails the selection of a subset of data from a larger dataset in order to draw inferences and make generalizations about the entire population. The utilization of this approach effectively mitigates the computing burden, particularly in scenarios involving extensive datasets such as social media. However, it is important to acknowledge that this strategy may introduce biases. Sampling can lead to biased insights due to the overrepresentation of certain user groups, opinions, or subjects [4]. These predispositions can have substantial impacts, particularly when these insights are employed to make critical decisions in domains such as marketing or policy formulation.

## 3. Methodology

1)Information Collection Process: In this research study, data was obtained from two prominent social networking platforms, namely Facebook and Twitter. The selection of these platforms was based on their active user participation and the diverse range of perspectives they facilitate [5]. To ensure contemporaneity in the analysis, a window spanning from January 1 to December 31, 2022, was established. During this time frame, a preliminary dataset consisting of 10 million posts and tweets was collected.

2)Belief Analysis Techniques: Two primary methodologies were employed for sentiment analysis: The first technique is the Lexicon-based Technique. The methodology employed in this approach involves quantifying the presence of positive and negative terms within a particular textual corpus, afterwards utilizing this information to calculate an overall assessment or rating. The SentiWordNet vocabulary, which is widely acknowledged, was utilized for this purpose [6]. The second technique is Artificial Intelligence. The designer of the study utilized a curated dataset consisting of 100,000 social media posts to train the model. This dataset included postings with positive, negative, and neutral sentiments. The idea was implemented using the Scikit-learn module in Python, making use of the Logistic Regression algorithm for its effectiveness in handling datasets with high dimensions.

3)Random Sampling Methodology: A complete analysis was conducted on a random sample of 1 million posts and tweets selected from the preliminary dataset. The selection process was conducted using the random module in Python, ensuring an equal probability of selection for each data point. In order to assure representation from both Twitter and Facebook, the sample procedure employed in the study was stratified by platform. This technique aims to depict the wider belief landscape while addressing computational limitations and ensuring quick analysis [7].

## 4. Types of Bias in Random Testing for Sentiment Analysis

Random sampling is a strategy that is widely employed to alleviate computing burdens. However, it is important to acknowledge that this approach may introduce various biases when conducting sentiment

analysis, particularly in the context of analyzing social media data. Understanding these predispositions is crucial for appropriately interpreting outcomes.

Firstly, let us consider the concept of choice bias. This phenomenon occurs when certain subsets of the data exhibit different odds of being selected for inclusion in the sample compared to others [8]. The possible bias in sentiment analysis arises from the overrepresentation of tweets including popular hashtags, which may be preferentially retrieved due to their higher exposure. Consider the impact of a popular hashtag, #AmazingProduct, which predominantly elicits positive sentiments. The act of excessively utilizing this particular hashtag has the potential to create an exaggerated perception of favorable sentiment towards a product.

Furthermore, the concept of protection bias should be considered. Protection bias is a phenomenon that arises when certain segments of the population are underrepresented or entirely absent from the sample [9]. Within the domain of social media, there exists a disparity in the frequency of material contribution across users. For instance, it is seen that younger users tend to exhibit a higher frequency of publishing content related to hot topics, whereas elderly users may display a somewhat lower inclination towards such activities. If the sampling technique fails to consider the difference in publication frequency, the viewpoints of younger users can dominate the results, perhaps leading to a biased comprehension of the overall attitude.

The third issue to be addressed is non-response. The term "predisposition" refers to an individual's inherent inclination or susceptibility towards a The emergence of this type of predisposition occurs when some segments of the population exhibit non-response or are less likely to be included in the sampling process [10]. Within the above context, consider a hypothetical scenario in which users or administrators of a platform engage in the act of removing negative evaluations or comments. If the dataset lacks these posts, the subsequent analysis may erroneously conclude a higher level of positive sentiment than what truly exists.

On the context of sentiment analysis on social media, it is crucial to acknowledge and mitigate biases due to the dynamic and diverse nature of emotions and opinions. This ensures that the insights obtained are a more accurate representation of the underlying landscape of beliefs.

## 5. Case Research Study: Evaluating Bias in Belief Analysis Outcomes

In order to elucidate the ramifications of biases in random sampling, a pragmatic case study was undertaken. The primary objective was to compare the findings of belief analysis obtained from a random sample with the outcomes derived from the entire dataset.

The comprehensive dataset consisted of 10 million posts and tweets collected from Twitter and Facebook during the year 2022. In accordance with the previously outlined technique, a random sample procedure was employed to choose a total of one million posts. The belief analysis was conducted on both datasets using identical machine learning and lexicon-based methodologies.

The study reveals that, with regards to General Sentiment Circulation, the overall dataset had beliefs distributed as follows: 40% positive, 35% neutral, and 25% unfavorable. Conversely, the random sample exhibited a distribution of 45% favorable, 30% neutral, and 25% unfavorable. This observation implies a potential overabundance of positive beliefs within the randomly selected sample.

Regarding Topic-Specific Beliefs, Upon isolating attitudes pertaining to the trending topic of "Social Media Personal Privacy," an examination of the complete dataset indicated a predominance of negative sentiments (60% unfavorable, 20% neutral, 20% favorable). However, the random sample exhibited a slight inclination towards negativity, with 50% of responses being negative, 25% neutral, and 25% positive [11]. This observation suggests the existence of choice bias, which is likely caused by the over-representation of tweets containing certain popular hashtags associated with positive feelings.

Regarding the phenomenon of Demographic Skewness, an analysis conducted on the posts based on age demographics unveiled a conspicuous lack of representation among individuals aged 50 years and older within the randomly selected sample. The opinions expressed by this particular cohort, which tended to be more neutral in comparison to younger age groups, were not sufficiently captured, suggesting a bias towards protection.

In conclusion, the utilization of random sampling in sentiment analysis has proven to be efficient, although it has also revealed significant biases inside the process. The results underscore the need of ensuring that the sample methodology is highly representative, particularly in situations when the potential consequences of misinterpretation are significant. Furthermore, it emphasizes the necessity of aligning any observations derived from experimental data with comprehensive studies of large datasets wherever feasible [12].

## 6. Reducing Bias in Random sampling: Techniques

It is imperative to address the biases that may arise from random sampling in order to ensure the accuracy and representativeness of belief analysis results. Various strategies can be employed to mitigate these predispositions.

Firstly, let us consider the concept of weighting methods. When there is an imbalance in the representation of various groups within a sample, the technique of weighting can be employed to modify the impact of particular data points on the ultimate analysis [13]. For instance, in the case when younger users are disproportionately represented within a randomly selected sample, it is possible to assign a reduced weight to each of their posts in order to mitigate the impact of their outlier presence. In an alternative approach, greater emphasis can be placed on under-represented groups by assigning them higher weights, hence amplifying their influence. This adjustment aims to align the sample more accurately with the actual distribution within the population.

Additionally, the implementation of stratified testing is worth considering. Stratified sampling is the process of splitting a population into subgroups, known as strata, that exhibit homogeneity. Subsequently, a random sample is selected from each stratum in proportion to its size [14]. When examining belief analysis on social networks, it is possible to construct strata by considering variables like as age, posting frequency, and platform utilization. By ensuring enough representation of each stratum, this strategy can provide a more equitable and accurate depiction of opinions across diverse populations.

Thirdly, the techniques of oversampling and undersampling. These techniques modify the sample's structure in order to address any imbalances. The process of oversampling entails increasing the number of samples obtained from groups that are under-represented, whereas undersampling involves reducing the number of samples obtained from groups that are over-represented [3]. In the event that posts conveying impartial beliefs are few yet significant, it may be necessary to oversample them in order to ensure their adequate representation in the ultimate analysis.

Furthermore, the use for enhanced precision is discussed. By employing these methodologies, researchers can significantly enhance the representativeness of their samples, hence yielding more accurate conclusions in sentiment analysis. An illustrative instance involves the utilization of stratified sampling, wherein beliefs are systematically captured across several demographic groups, hence enabling a more comprehensive depiction of public opinion. When employed with discretion, weighting procedures can rectify instances of over- or under-representations, ensuring that the final analysis accurately reflects the true distribution of beliefs within the broader dataset.

The incorporation of these strategies can have a significant impact on the differentiation between an analysis that yields authentic insights and one that may mislead due to inherent biases in the sampling process.

## 7. Discussion

### 7.1. Implications and Applications

The recognition of potential biases in random sampling, particularly in the context of belief analysis, is crucial for facilitating precise and significant analyses of data. The following are the broader implications and potential applications:

In the realm of corporate decision-making, it is common practice for organizations to rely on belief analysis as a means of assessing the public's perception and comprehension of their offerings, be it products, services, or branding initiatives. The presence of a consciousness regarding potential biases

serves to prevent corporations from misinterpreting the dominant belief, hence enabling more accurate adaptations to marketing strategies or product alterations.

Policy Solution: Public and private entities frequently employ sentiment analysis as a strategic tool to leverage public attitude towards policies or events. The acknowledgement and depiction of sample biases ensures that policies be formulated on the basis of a genuine comprehension of public opinion, rather than a distorted one.

In the realm of academic research, it is imperative for scholars and researchers to incorporate belief analysis in order to enhance the validity of their studies and derive more dependable findings. This can be achieved by taking into consideration any potential sample biases that may exist. This practice guarantees the reliability of scholarly contributions and enhances the possibility of practical implementations based on their discoveries.

Enhanced Machine Learning Designs: When conducting training on machine learning models using sentiment data, it is crucial to utilize a balanced dataset that is free from substantial biases. This approach ensures that the model will exhibit strong generalization capabilities when applied to novel and unseen data. Having an understanding of the potential bias introduced by random sampling can be beneficial at the information preprocessing stage, leading to the development of more resilient architectures.

Stakeholder Trust: In the context of firms or platforms that share belief analysis data, establishing transparency regarding the presence of sampling bias and the measures taken to address it helps foster trust among stakeholders. The demonstration of due care is evident in the comprehensive and inclusive nature of the analysis.

In summary, understanding the intricacies of random sample bias not only enhances the precision of sentiment analysis outcomes but also strengthens the decision-making process across many industries. In an era characterized by the significance of data-driven insights, the certainty of their foundation in rigorous and impartial research holds immense value.

### 7.2. Future Instructions and Difficulties

The domain of sentiment analysis, namely within the realm of social network data, is rapidly advancing. The aforementioned progress presents intriguing opportunities for further investigation, as well as ongoing challenges that necessitate careful consideration.

*7.2.1. Future Instructions.* First, Deep Learning and Predisposition Detection. As the complexity of deep learning models increases, there is a potential to develop architectures that possess the ability to promptly detect and adapt to biases present in sentiment analysis datasets, so ensuring more resilient and reliable results.

Second, Cross-Platform Analysis. Given the diverse range of social networking platforms available, each characterized by distinct user behavior and demographics, it is imperative to do study on the manifestation of belief predispositions across these platforms. Furthermore, it is essential to explore methods of harmonizing these predispositions to facilitate a more coherent analysis of beliefs.

Temporal Predisposition Analysis: Given the volatile and dynamic nature of beliefs expressed on social media, it would be valuable for future research to investigate the influence of temporal biases arising from the timing of data collection on the conclusions of belief analysis.

Integrating Contextual Details: The utilization of metadata and contextual information, such as geography, user demographics, or neighboring posts, can provide more comprehensive insights and aid in comprehending the nuances of sentiment predispositions. The investigation of strategies for seamlessly integrating this data may be of utmost importance.

*7.2.2. Difficulties.* One of the primary challenges encountered is the sheer volume and velocity of information. The acquisition of representative samples and the maintenance of up-to-date belief assessments face challenges due to the substantial amount of data available on social network platforms and the rapid pace at which new data is generated.

The phenomenon of multilingualism and cultural subtleties has a significant role in the context of social networks, as these platforms have a global reach and feeling expressed within them can be strongly influenced by cultural and linguistic factors. The enormous challenge of overcoming biases that arise from language translation or cultural misinterpretations persists.

The following aspect pertains to the employment of sarcasm and subtlety. The identification of nuanced emotions, such as sarcasm or irony, remains a significant challenge, even when employing advanced algorithms. The potential for biases arises when there is a possibility of misinterpretation or misclassification of these thoughts.

The remaining category pertains to Platform Algorithms. Social networking platforms employ algorithms to decide the extent to which content is exposed. These algorithms have the potential to bias the available data for sampling, resulting in platform-induced predispositions.

The final aspect to be taken into account is the ethical considerations. As researchers strive to obtain more inclusive samples, they must confront ethical challenges, such as privacy considerations and ensuring that data gathering procedures are clear and respectful of user rights.

In conclusion, although significant progress has been achieved in the domain of sentiment analysis, the dynamic characteristics of social media and the complex structure of human beliefs ensure that the investigation of random sample bias will remain a pertinent subject of scholarly inquiry. The future holds potential for advancements that can enhance our comprehension and evaluation of sentiments. However, it is crucial for academics and professionals to remain vigilant regarding the biases that may arise and potentially distort these findings.

## 8. Conclusion

In the contemporary era of digitalization, whereby decision-making processes across several domains such as marketing and public law heavily rely on data, sentiment analysis emerges as a promising avenue for extracting valuable insights. Extracting the collective sentiments derived from extensive repositories of social network data yields unparalleled insights into prevailing public attitude, consumer inclinations, and nascent trends. However, this study highlights that the techniques employed to derive these insights are susceptible to potential errors, with random sample bias being a prominent concern.

The key findings derived from our investigation can be summarized into four main themes. The first concept under consideration is the Universality of Predisposition. The issue of random sampling propensity is not an occasional or peripheral matter. The pervasive nature of its existence, if left unchecked, has the potential to significantly distort our comprehension of sentiments expressed on social media sites. The second aspect is to the diversity of bias types. Various factors might contribute to misleading belief results, ranging from individual predispositions towards choosing or non-response. Recognizing and understanding these categories constitutes the initial phase in the process of mitigation. The third point asserts that mitigation is a viable option. Despite the presence of various challenges, there are well-established methodologies, such as weighting techniques and stratified sampling, that can be utilized to mitigate the impact of bias and enhance the accuracy of sentiment analysis, hence facilitating a more accurate depiction of public opinion. The final aspect to consider is the real-world implications. The potential misunderstandings arising from biased sentiment analysis can result in tangible repercussions. The consequences are consistently significant, whether it is a situation where a company misjudges the response to a newly introduced product or legislators fail to accurately gauge public opinion.

In the future, as the realm of social media continues to evolve and grow, the importance of rigor in the examination of beliefs will inevitably increase. The present study functions as a dual-purpose narrative, serving as a cautionary anecdote while also providing guidance, so emphasizing the importance of methodological rigor. In order to obtain truly important insights, it is imperative that our assessments not only possess a profound level of depth, but also rest upon a firm foundation of objective knowledge. The efficacy of sentiment analysis is not solely derived from the act of monitoring online individuals, but rather from the ability to genuinely understand them. As previously emphasized, a

comprehensive comprehension can only be attained through the analysis of unbiased material, devoid of any potential biases stemming from random sampling.

## References

[1]  Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. Foundations and Trends® in Information Retrieval, 2(1–2), 1–135.

[2]  Kemp, S. (2020). Digital 2020: Global digital overview. Datareportal.

[3]  Agarwal, A., Xie, B., Vovsha, I., Rambow, O., & Passonneau, R. (2011). Sentiment analysis of Twitter data. In Proceedings of the workshop on languages in social media (pp. 30-38). Association for Computational Linguistics.

[4]  Morstatter, F., Pfeffer, J., Liu, H., & Carley, K. M. (2013). Is the sample good enough? Comparing data from Twitter's streaming API with Twitter's firehose. In Seventh international AAAI conference on weblogs and social media.

[5]  Zhang, A. X., Chen, R. M., & Carley, K. M. (2018). Large Scale Structure and Dynamics of Complex Networks: From Information Technology to Finance and Natural Science. World Scientific.

[6]  Baccianella, S., Esuli, A., & Sebastiani, F. (2010). SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10). European Language Resources Association (ELRA).

[7]  Smith, T. M. (2013). Sampling and statistical methods for behavioral ecologists. Cambridge University Press.

[8]  Bethlehem, J. (2010). Selection bias in web surveys. International Statistical Review, 78(2), 161-188.

[9]  Wang, W., Rothschild, D., Goel, S., & Gelman, A. (2015). Forecasting elections with non-representative polls. International Journal of Forecasting, 31(3), 980-991.

[10]  Groves, R. M. (2006). Nonresponse rates and nonresponse bias in household surveys. Public Opinion Quarterly, 70(5), 646-675.

[11]  Salganik, M. J. (2017). Bit by bit: Social research in the digital age. Princeton University Press.

[12]  Lohr, S. (2019). Sampling: Design and Analysis. Chapman and Hall/CRC.

[13]  Cochran, W. G. (2007). Sampling techniques. John Wiley & Sons.

[14]  Chawla, N. V. (2005). Data mining for imbalanced datasets: An overview. In Data mining and knowledge discovery handbook (pp. 853-867). Springer, Boston, MA.