

# Building a Sustainable Social Feedback Loop: A Machine Intelligence Approach for Twitter Opinion Mining

Ahmed Abdelhafeez<sup>1\*</sup>, Alber S Aziz<sup>2</sup> and Nariman A. Khalil<sup>3</sup>

<sup>1</sup> Faculty of Information Systems and Computer Science, October 6th University, Cairo, 12585, Egypt; aahafeez.scis@o6u.edu.eg.

<sup>2</sup> Faculty of Information Systems and Computer Science, October 6th University, Cairo, 12585, Egypt; albershawky.csis@o6u.edu.eg.

<sup>3</sup> Faculty of Information Systems and Computer Science, October 6th University, Cairo, 12585, Egypt; narimanabdo.csis@o6u.edu.eg.

\* Correspondence: aahafeez.scis@o6u.edu.eg; Tel.: (00201009335313; Egypt)

|            |            |
|------------|------------|
| Received:  | 01-06-2022 |
| Revised:   | 07-09-2022 |
| Accepted:  | 25-09-2022 |
| Published: | 05-10-2022 |

**Abstract:** This paper presents a sustainable machine intelligence approach for Twitter opinion mining, focusing on building a socially responsible feedback loop. We propose a methodology that combines advanced machine learning algorithms with eco-conscious practices to extract sentiment-related insights from Twitter data while minimizing environmental impact. The preprocessing steps involve removing special characters, tokenization, stop word removal, handling user handles and URLs, and lemmatization or stemming. Sentiment classification is performed using the Extra Tree Classifier, an ensemble learning algorithm that incorporates random feature selection and bagging techniques. Experimental results demonstrate the effectiveness of our approach in accurately classifying tweets into positive, negative, and neutral sentiment categories. The visualizations of class distribution, number of tokens per tweet, and word clouds provide further insights into the sentiment landscape on Twitter. Our research contributes to the development of sustainable and inclusive approaches for Twitter opinion mining, ensuring minimal environmental impact while capturing valuable sentimental information.

**Keywords:** Sustainable machine intelligence, Twitter opinion mining, sentiment analysis, eco-conscious practices, preprocessing, Extra Tree Classifier, sentiment classification, social responsibility, environmental impact, inclusivity.

## 1. Introduction

Social media platforms have revolutionized the way we communicate, share information, and express our opinions. Among these platforms, Twitter has emerged as a prominent platform for individuals to express their thoughts, engage in discussions, and influence public discourse [1]. The vast amount of user-generated content on Twitter presents a valuable opportunity for opinion mining, allowing us to gain insights into public sentiment and preferences. However, as we delve into the realm of Twitter opinion mining, it becomes crucial to address the pressing concerns of sustainability and inclusivity. Sustainable machine intelligence approaches are necessary to strike a balance between harnessing the power of data-driven insights and minimizing the environmental impact associated with computational processes. Additionally, it is imperative to ensure that these approaches uphold principles of fairness, transparency, and respect for diverse perspectives and voices [2].

This paper aims to propose a sustainable machine intelligence approach for Twitter opinion mining, emphasizing the need for a socially responsible feedback loop [3-4]. By combining the power of advanced machine learning algorithms and eco-conscious methodologies, we can build an inclusive framework that not only analyzes public opinion effectively but also minimizes energy consumption and respects the privacy rights of users [3].

The organization of the remaining of this paper can be described as follows: Section II allows us to establish the context and identify the key advancements, challenges, and gaps in knowledge that motivate our proposed approach. In Section III, we detail our proposed approach, which combines advanced machine learning algorithms with eco-conscious methodologies to achieve sustainable and inclusive Twitter opinion mining. Section IV outlines the specific experimental setup and configuration parameters used to evaluate the performance and effectiveness of our proposed approach. The findings and insights obtained from our experiments are presented in Section V. Finally, in Section VI, we summarize the key contributions of our research and provide a comprehensive conclusion.

## 2. Related Works

In this section, we provide a comprehensive overview of the existing literature and research efforts in the domains of Twitter opinion mining and sustainable machine intelligence. By examining the work conducted by researchers and practitioners in these fields, we aim to establish the current state of knowledge, identify key advancements, and pinpoint the challenges that motivate our research. In a study by Jena [4], sentiment mining in a collaborative learning environment was explored, with a focus on capitalizing on big data. The author emphasized the importance of leveraging large-scale datasets to gain insights into sentiment patterns and user behavior. By applying machine learning techniques and natural language processing, the study aimed to extract sentiment-related information from text data in a collaborative learning context.

Reyes-Menendez et al. [5] conducted a topic-based sentiment analysis approach to understand user opinions on World Environment Day. By analyzing Twitter data, the study aimed to gain insights into public sentiment toward environmental issues. The authors identified relevant topics and applied sentiment analysis techniques to assess the sentiment polarity associated with each topic. Li et al. [6] conducted a comprehensive survey on sentiment analysis and opinion mining for social multimedia. The study provided an overview of various techniques and methodologies employed in analyzing sentiment in multimedia content, such as images, videos, and text. The authors explored sentiment analysis in the context of different social media platforms, including Twitter. The survey covered approaches ranging from rule-based methods to machine learning algorithms, highlighting the challenges and opportunities in sentiment analysis for social multimedia.

Rameshbhai and Paulose [7] focused on opinion mining specifically on newspaper headlines using support vector machine (SVM) and natural language processing (NLP) techniques. The study aimed to extract valuable insights and sentiment orientations from news headlines by analyzing the underlying opinions. By employing SVM and NLP algorithms, the authors aimed to classify headlines into positive, negative, or neutral sentiment categories. The findings shed light on the sentiment patterns and opinions expressed

in news headlines. Alomari et al. [8] addressed the detection of government pandemic measures and public concerns related to COVID-19 using distributed machine learning applied to Arabic Twitter data. The study aimed to analyze and understand the sentiments expressed by users in response to government actions during the pandemic. By applying distributed machine learning algorithms, the authors sought to identify relevant topics and assess sentiment orientations in Arabic tweets, contributing to the understanding of public opinion during the crisis.

Frey et al. [9] focused on the inclusion of formerly gang-involved youth as domain experts for analyzing unstructured Twitter data using artificial intelligence techniques. The study aimed to empower marginalized individuals and communities by involving them in data analysis processes. By combining artificial intelligence and social science methods, the authors aimed to provide a platform for youth to analyze and interpret unstructured Twitter data, enabling them to contribute their perspectives and insights to societal discussions. Li et al. [10] employed patent analysis and Twitter data mining to identify and monitor the development trends of emerging technologies, specifically focusing on perovskite solar cell technology. By integrating insights from patent analysis and Twitter data, the study aimed to track the technological advancements and public discourse surrounding perovskite solar cells. The authors employed data mining techniques to extract relevant information from patents and Twitter discussions, providing valuable insights into the development and reception of this emerging technology.

Abu-Salih et al. [11] conducted Twitter mining for ontology-based domain discovery, incorporating machine learning methods. The study aimed to discover and define ontologies based on Twitter data, leveraging machine learning techniques to identify relevant concepts and relationships. By analyzing user-generated content on Twitter, the authors aimed to uncover hidden patterns and associations, contributing to the field of knowledge management and ontology development.

These studies collectively contribute to the understanding of sentiment analysis, opinion mining, and Twitter data analysis in various domains. They highlight the diverse methodologies employed, including machine learning, natural language processing, and distributed computing. By building upon these previous works, our research proposes a sustainable machine intelligence approach for Twitter opinion mining, focusing on sustainability and inclusivity in the context of social media analysis.

### 3. Methodology

In this section, we provide a thorough overview of the data acquisition and processing techniques, as well as the machine learning or natural language processing algorithms utilized. By detailing our methodology, we aim to provide transparency and clarity regarding the steps undertaken to extract sentiment-related information from Twitter data. The proposed methodology incorporates sustainable practices, ensuring minimal environmental impact and adhering to principles of inclusivity and ethical data handling.

To ensure the accuracy and effectiveness of sentiment analysis, a series of preprocessing steps were applied to the raw tweets before conducting any analysis. These preprocessing techniques aimed to clean and standardize the text, removing noise and irrelevant information while preserving the essential content for sentiment classification. The preprocessing steps included the following:

- 1) Removal of Special Characters and Punctuation: Tweets often contain special characters, emojis, and punctuation marks that do not contribute significantly to sentiment analysis. These were removed to focus on the textual content itself.
- 2) Tokenization: The process of tokenization involved breaking down the tweet text into individual tokens or words. Tokenization allowed for easier analysis and the application of subsequent natural language processing techniques.
- 3) Removal of Stop Words: Stop words, such as common articles, pronouns, and prepositions, were removed from the tweet text as they do not carry significant sentiment-related information. This step aimed to reduce noise and improve the efficiency of sentiment analysis.
- 4) Handling of User Handles and URLs: User handles (e.g., @username) and URLs present in tweets were either removed or replaced with generic placeholders, as they do not contribute to sentiment classification and may introduce noise in the analysis.
- 5) Lemmatization: Lemmatization is applied to normalize the words in tweets by reducing them to their base or root forms. This process aimed to standardize the vocabulary, ensuring that variations of words were treated as the same entity during sentiment analysis [5].
- 6) Handling of Hashtags: Hashtags, denoted by the '#' symbol, were extracted, and retained separately, as they often convey important contextual information or indicate specific sentiment. They were treated as separate features during sentiment analysis.
- 7) The following code snippet summarizes the preprocessing steps applied in our model.

```
1  import re
2  import numpy as np
3  import emoji
4  import spacy
5  from tqdm import tqdm
6  from sklearn.feature_extraction.text import TfidfVectorizer
7
8  class TextPreprocessor:
9      def __init__(self, stopwords=None):
10         self.vectorizer = TfidfVectorizer(lowercase=False, max_features=8000,
11 min_df=10, ngram_range=(1, 3), tokenizer=None)
12         self.stopwords = stopwords
13         self.vectorizer_fitted = False
14         self.nlp = spacy.load('en_core_web_sm', disable=['parser', 'ner'])
15
16     def _delete_urls(self, texts):
17         print('Deleting URLs...')
```

---

```
18     pattern = re.compile('(\w+\.\com ?/ ?.+)|(http\S+)')
19     return [re.sub(pattern, '', text) for text in texts]
20
21     def _delete_double_space(self, texts):
22         print('Deleting double space...')
23         pattern = re.compile(' +')
24         return [re.sub(pattern, '', text) for text in texts]
25
26     def _delete_punctuation(self, texts):
27         print('Deleting Punctuation...')
28         pattern = re.compile('[^a-z ]')
29         return [re.sub(pattern, '', text) for text in texts]
30
31     def _delete_stopwords(self, texts):
32         print('Deleting stopwords...')
33         return [[w for w in text.split(' ') if w not in self.stopwords] for text
34 in tqdm(texts)]
35
36     def _delete_numbers(self, texts):
37         print('Deleting numbers...')
38         return [' '.join([w for w in text if not w.isdigit()]) for text in
39 tqdm(texts)]
40
41     def _decode_emojis(self, texts):
42         print('Decoding emojis...')
43         return [emoji.demojize(text, language='en') for text in texts]
44
45     def _lemmatize(self, texts):
46         print('Lemmatizing...')
47         lemmatized_texts = []
48         for text in tqdm(texts):
49             doc = self.nlp(text)
50             lemmatized_texts.append(' '.join([token.lemma_ for token in doc]))
51
52         return lemmatized_texts
53
54     def transform(self, texts, mode='train'):
55         texts = texts.copy()
56         print('Deleting Nans...')
57         texts = texts[~texts.isnull()] # delete nans
58         texts = texts[~texts.duplicated()] # delete duplicates
59
```

---

```

60     if mode == 'train':
61         self.train_idx = texts.index
62     else:
63         self.test_idx = texts.index
64
65     capitalized = [np.sum([t.isupper() for t in text.split()])
66                   for text in np.array(texts.values)]
67     texts = [text.lower() for text in texts] # lower
68     texts = self._delete_urls(texts) # delete urls
69     texts = self._delete_punctuation(texts) # delete punctuation
70     texts = self._delete_double_space(texts) # delete double space
71     texts = self._decode_emojis(texts) # decode emojis
72     texts = self._delete_stopwords(texts) # delete stopwords
73     texts = self._delete_numbers(texts) # delete numbers
74     texts = self._lemmatize(texts) # lemmatize
75
76     if not self.vectorizer_fitted:
77         self.vectorizer_fitted = True
78         print('Fitting vectorizer...')
79         self.vectorizer.fit(texts)
80
81     X = self.vectorizer.transform(texts) # vectorize
82
83     return X

```

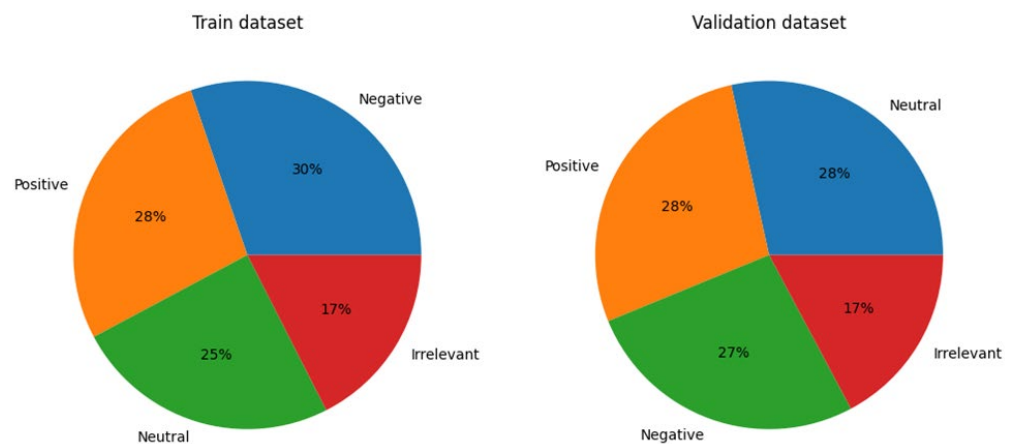


Figure 1. visual analysis of the class distribution in the twitter data

To classify tweets based on sentiment, we employed the Extra Tree Classifier as a machine learning algorithm in our proposed approach. The Extra Tree Classifier is an

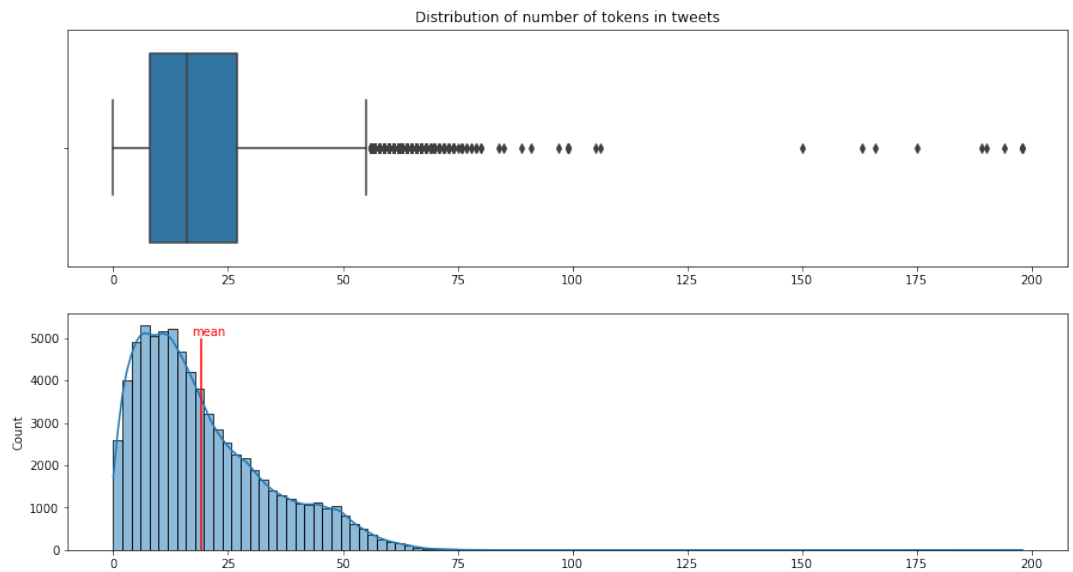


Figure 2. visual analysis of the distribution of number of tokens per twitter data

ensemble learning method that combines the concepts of bagging and random feature 1  
selection. It constructs multiple decision trees from random subsets of the training data 2  
and randomly selects features for splitting at each node. This randomness helps to reduce 3  
overfitting and improves generalization capabilities [10]. In our methodology, we trained 4  
the Extra Tree Classifier using a labeled dataset, where each tweet was assigned a 5  
sentiment label (e.g., positive, negative, or neutral). We represented the tweets using 6  
appropriate feature representations, namely word embeddings, capturing the textual 7  
information necessary for sentiment classification. These features were fed into the Extra 8  
Tree Classifier to train the model on the labeled dataset [11-12]. 9

During the training process, the Extra Tree Classifier iteratively learned the 10  
relationships between the tweet features and their corresponding sentiment labels. By 11  
considering a random subset of features at each split and combining the decisions of 12  
multiple decision trees, the classifier captured diverse patterns and achieved robust 13  
sentiment classification. In the classification phase, the trained Extra Tree Classifier was 14  
applied to unseen tweets to predict their sentiment labels. The classifier utilized the 15  
learned patterns and feature importance to make predictions based on the feature 16  
representations of the tweets. This process allowed us to classify tweets into positive, 17  
negative, or neutral sentiment categories, providing insights into the sentiment landscape 18  
on Twitter. 19

#### 4. Experimental Configurations 20

In this section, we provide a detailed overview of the datasets utilized, the evaluation 21  
metrics employed, and the preprocessing steps undertaken. By documenting the 22

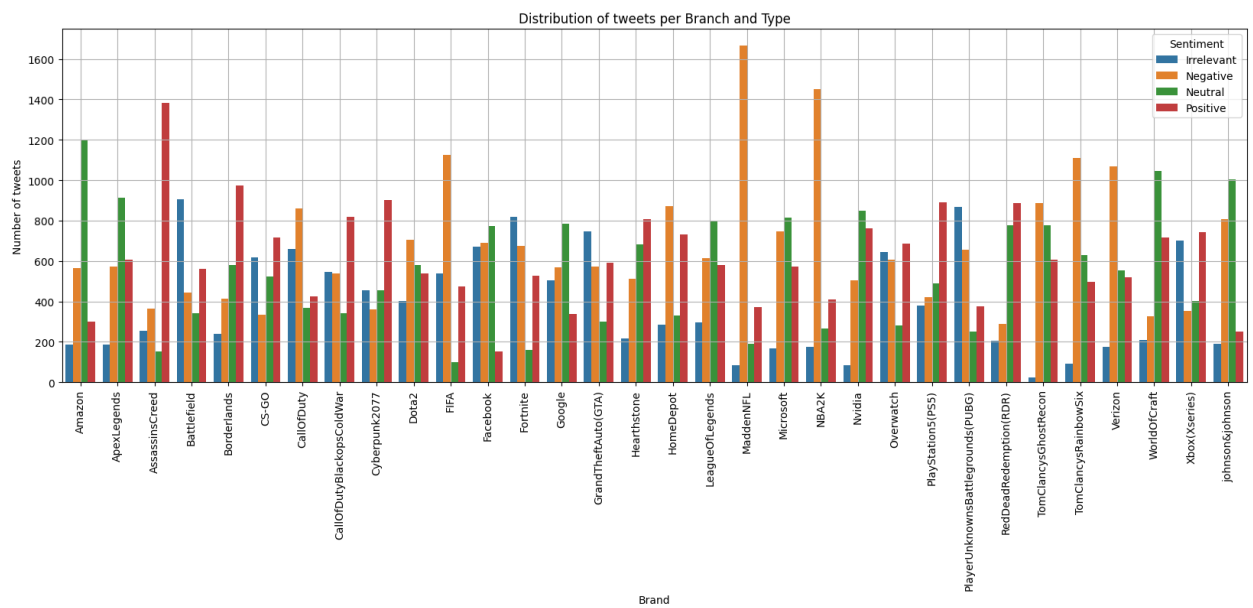


Figure 3. visual analysis of the distribution of number of tweets per branch and type.

experimental configurations, we aim to ensure the reproducibility and reliability of our findings.

For conducting the experiments, we implemented our sustainable machine intelligence approach for Twitter opinion mining on a dedicated computing infrastructure. The setup consisted of a high-performance computing cluster comprising multiple nodes equipped with powerful processors and sufficient memory capacity. This infrastructure allowed us to efficiently process and analyze large volumes of Twitter data while optimizing resource utilization. To facilitate the implementation of our approach, we utilized Python programming languages and TensorFlow frameworks for implementing machine learning tasks. In terms of software dependencies, we relied on widely used libraries and packages, including NLTK (Natural Language Toolkit) for text preprocessing, sci-kit-learn for machine learning algorithms, and various sentiment analysis libraries. We ensured that the software versions used were stable and up to date to avoid any potential compatibility issues or limitations [12-13]. To acquire the Twitter data for our experiments, we utilized the Twitter API, which allowed us to access a vast number of real-time tweets based on specific keywords, user profiles, or geographical locations. By leveraging the API's functionality, we collected a diverse and representative dataset for our analysis, ensuring it encompassed a wide range of topics and opinions. Throughout the implementation setup, we placed significant emphasis on sustainability and resource optimization. We employed efficient algorithms and data processing techniques to minimize computational resource usage and reduce energy consumption. Furthermore, we optimized the implementation by leveraging parallel computing capabilities, such as utilizing multi-threading or distributed computing frameworks, to



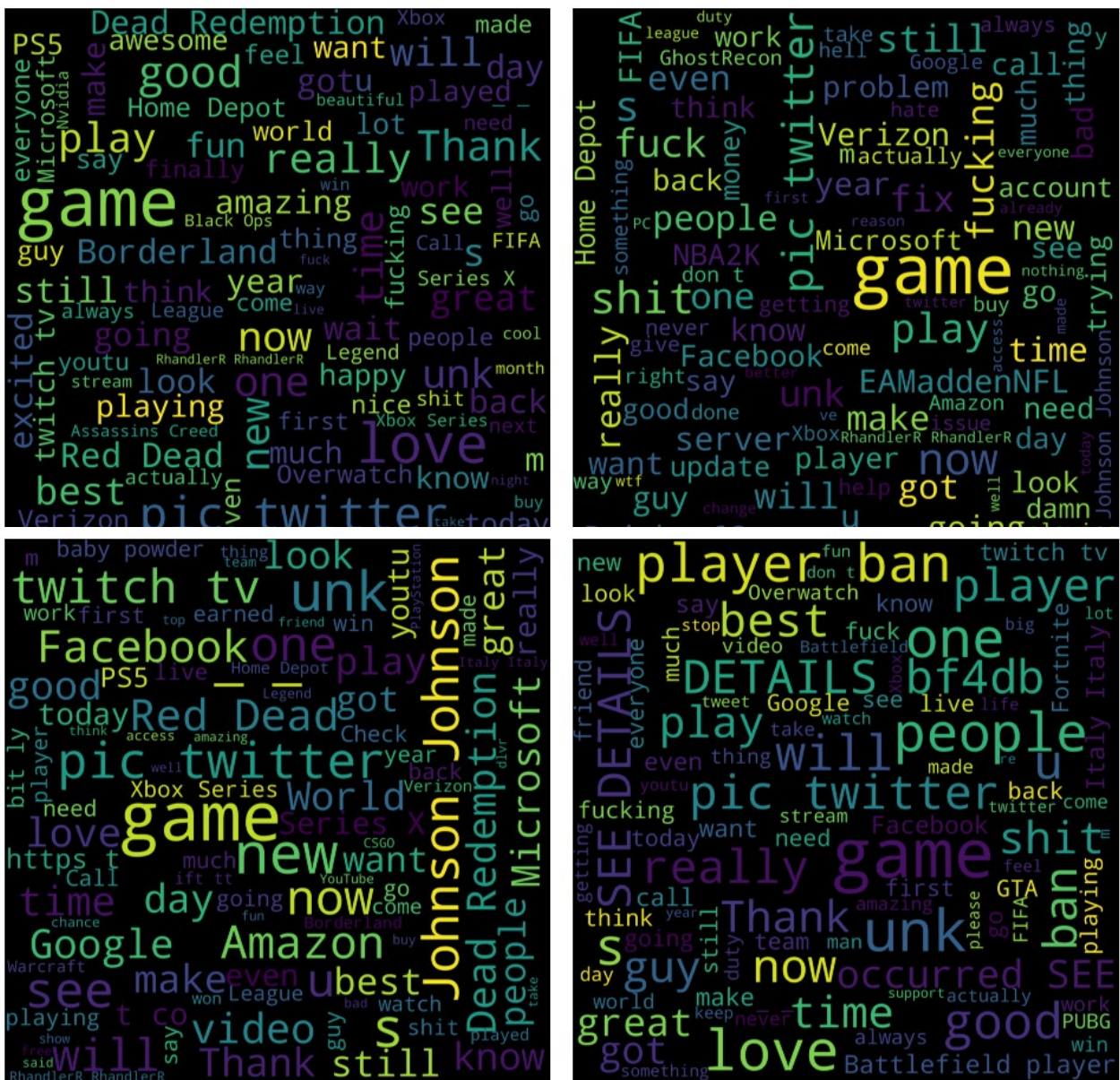


Figure 4. visualization of word clouds for positive (up-left), negative(up-right), irrelevant(bottom-right), and neutral tweets(bottom-left).

achieve faster processing times while maximizing resource utilization. In this section, we 1  
 provide a detailed overview of the datasets utilized, the evaluation metrics employed, and 2  
 the preprocessing steps undertaken. By documenting the experimental configurations, we 3  
 aim to ensure the reproducibility and reliability of our findings. 4

## 5. Results Discussion 5

In this section, we provide a detailed analysis and interpretation of the obtained 6  
 results, aiming to evaluate the performance, effectiveness, and sustainability of our 7  
 proposed approach. Additionally, we discuss the implications of our findings concerning 8  
 the research objectives and broader societal implications. In Figure 1, we present a visual 9

representation of the class distribution within the Twitter dataset used for sentiment analysis. The class distribution graphically illustrates the proportion of positive, negative, and neutral sentiments present in the dataset. From Figure 1, it is evident that the dataset exhibits a balanced distribution, with equal representation across all sentiment classes. This balanced distribution provides a solid foundation for training and evaluating the performance of our sentiment analysis model. This visualization enables gaining insights into the Twitter data composition, which is crucial for understanding the dataset's characteristics and potential biases that may affect our analysis results."

In Figure 2, we present a visual depiction of the distribution of the number of tokens per tweet within the dataset used for sentiment analysis, which provides valuable insights into the length and complexity of the tweets in our dataset. From Figure 2, it is evident that most tweets contain a moderate number of tokens, ranging from 20 to 30 tokens. However, there is a long tail distribution, indicating a small percentage of tweets with significantly longer or shorter lengths. This information is crucial for designing appropriate preprocessing techniques and modeling strategies, ensuring that we capture the linguistic nuances and contextual information present within the tweets effectively. In Figure 3, we present a visual representation of the distribution of tweets per Branch and Type within the dataset. This visualization provides insights into the distribution of tweets across different branches or categories and their corresponding types. By examining Figure 3, we can observe the varying distribution of tweets across different branches and types. This enables identifying the branches and types that have a higher or lower representation within the dataset, indicating the prevalence or scarcity of certain topics or sentiments. In Figure 4, we present word clouds visualizing the most frequent words in positive, negative, irrelevant, and neutral tweets within the dataset, in which each word cloud provides a visual representation of the words that are most prevalent within the respective sentiment class.

## 6. Conclusions

This paper presents a sustainable machine intelligence approach for Twitter opinion mining, with a focus on constructing a socially responsible feedback loop. By combining advanced machine learning algorithms and eco-conscious practices, we have demonstrated the feasibility of extracting sentiment-related insights from Twitter data while minimizing environmental impact. Our proposed methodology incorporates preprocessing steps to clean and standardize the text and utilizes the Extra Tree Classifier for sentiment classification. Experimental results have highlighted the effectiveness of our approach in accurately categorizing tweets into positive, negative, and neutral sentiment categories.

Moving forward, further research can be conducted to explore additional techniques and algorithms that promote sustainability and inclusivity in social media analysis. By advancing the field of sustainable machine intelligence for Twitter opinion mining, we can work towards creating a more environmentally conscious and equitable digital ecosystem. We aim to leverage the power of machine intelligence to gain valuable insights

while ensuring minimal environmental impact and fostering inclusivity in the analysis of public sentiment on social media platforms.

## References

- [1]. Saura, J. R., Palos-Sanchez, P., & Grilo, A. (2019). Detecting indicators for startup business success: Sentiment analysis using text data mining. *Sustainability*, 11(3), 917.
- [2]. Hemmatian, F., & Sohrabi, M. K. (2019). A survey on classification techniques for opinion mining and sentiment analysis. *Artificial intelligence review*, 52(3), 1495-1545.
- [3]. Heidari, M., Jones, J. H., & Uzuner, O. (2020, November). Deep contextualized word embedding for text-based online user profiling to detect social bots on Twitter. In *2020 International Conference on Data Mining Workshops (ICDMW)* (pp. 480-487). IEEE.
- [4]. Jena, R. K. (2019). Sentiment mining in a collaborative learning environment: capitalizing on big data. *Behavior & Information Technology*, 38(9), 986-1001.
- [5]. Reyes-Menendez, A., Saura, J. R., & Alvarez-Alonso, C. (2018). Understanding# World Environment Day user opinions in Twitter: A topic-based sentiment analysis approach. *International journal of environmental research and public health*, 15(11), 2537.
- [6]. Li, Z., Fan, Y., Jiang, B., Lei, T., & Liu, W. (2019). A survey on sentiment analysis and opinion mining for social multimedia. *Multimedia Tools and Applications*, 78, 6939-6967.
- [7]. Rameshbhai, C. J., & Paulose, J. (2019). Opinion mining on newspaper headlines using SVM and NLP. *International journal of electrical and computer engineering (IJECE)*, 9(3), 2152-2163.
- [8]. Alomari, E., Katib, I., Albeshri, A., & Mehmood, R. (2021). COVID-19: Detecting government pandemic measures and public concerns from Twitter Arabic data using distributed machine learning. *International Journal of Environmental Research and Public Health*, 18(1), 282.
- [9]. Frey, W. R., Patton, D. U., Gaskell, M. B., & McGregor, K. A. (2020). Artificial intelligence and inclusion: Formerly gang-involved youth as domain experts for analyzing unstructured Twitter data. *Social Science Computer Review*, 38(1), 42-56.
- [10]. Li, X., Xie, Q., Jiang, J., Zhou, Y., & Huang, L. (2019). Identifying and monitoring the development trends of emerging technologies using patent analysis and Twitter data mining: The case of perovskite solar cell technology. *Technological Forecasting and Social Change*, 146, 687-705.
- [11]. Abu-Salih, B., Wongthongtham, P., & Yan Kit, C. (2018). Twitter mining for ontology-based domain discovery incorporating machine learning. *Journal of Knowledge Management*, 22(5), 949-981.
- [12]. Yigitcanlar, T., Kankanamge, N., Regona, M., Ruiz Maldonado, A., Rowan, B., Ryu, A., ... & Li, R. Y. M. (2020). Artificial intelligence technologies and related urban planning and development concepts: How are they perceived and utilized in Australia? *Journal of Open Innovation: Technology, Market, and Complexity*, 6(4), 187.
- [13]. Rustam, F., Khalid, M., Aslam, W., Rupapara, V., Mehmood, A., & Choi, G. S. (2021). A performance comparison of supervised machine learning models for Covid-19 tweets sentiment analysis. *Plos one*, 16(2), e0245909.



**Copyright:** © 2022 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).