

# Natural Language Processing and Text Analytics

## Uncovering Sentiment Patterns in ChatGPT

Exploring Twitter users' opinion about ChatGPT with sentiment analysis  
and topic modeling



Andrea Pérez López (anpe22aq) (158798)

Davide Maria Piva (dapi22ab) (158388)

Jenő Tóth (jeto22ac) (158386)

S M Ahasanul Karim (smka22ab) (158793)

Characters count: 26,912

Page count: 14

May 29, 2023

**Natural Language Processing Written Exam**

MSc. in Business Administration and Data Science

## Abstract

*This research concerns with finding out people’s reaction to the introduction of ChatGPT, an AI chatbot by OpenAI, by utilizing Natural Language Processing to examine Twitter posts related to the topic. The research aims to explore the sentiment and topics expressed in a dataset of more than sixty thousand tweets. Our research focuses on identifying the most common feelings towards ChatGPT through sentiment analysis and topic modeling. Our results provide valuable insights into public attitudes towards ChatGPT. Initial analysis revealed that 45% of users portrayed a positive attitude, while 41% remained neutral, and a minority of 14% demonstrated a negative perspective. Interestingly, the subjectivity analysis indicated that the Twitter discourse was mainly objective (64%), suggesting that users primarily share factual information about ChatGPT, rather than expressing personal sentiment or opinion. Analysis of word embeddings further supported these findings, demonstrating a general trend toward positive sentiment. The paper examines daily changes in attitude, and finds that positive sentiments outnumbered negative ones every day. Overall, our findings suggest a prevailing positive perception of ChatGPT among users, with a significant proportion of the discourse focusing on the tool’s factual aspects.*

**Keywords:** Sentiment Analysis, Natural Language Processing, Topic Modeling, ChatGPT, Twitter, Word2vec, Top2Vec

# 1 Introduction

The introduction of ChatGPT to the general public marks a milestone in the diffusion of Artificial Intelligence (AI) and can be considered one of the most significant achievements for AI tools. This chatbot relies on an advanced natural language model to interact with users, providing them with answers in a conversational manner.

For the first time, the accessibility and utilization of such a tool have made it evident the exceedingly rapid progress to which AI is subjected nowadays. This not only highlights the numerous implementations that humanity may witness in the years to come but also emphasizes the social, economic, and technical risks that may arise.

Therefore, this paper aims to understand people's attitudes towards ChatGPT by analyzing a dataset comprising 66,375 tweets from around the world (Domingo, 2023), focusing on this topic through Natural Language Processing (NLP) techniques. Specifically, sentiment analysis is implemented to evaluate the most common reactions and feelings that people have towards ChatGPT. In addition, topic modeling is employed to gain a broader insight into the dataset, identifying recurring patterns or ideas. These steps will be implemented using Python's Word2vec and Top2Vec packages.

## 2 Literature Review

The increasing availability of digital textual data has led to the widespread application of sentiment analysis across various domains. According to Mejova (2009), sentiment analysis is commonly regarded as a sub-field encompassing natural language processing, text mining, and computational linguistics. In their paper titled "Sentiment Analysis: An Overview," the author provides a comprehensive description of this technique, which aims to identify and extract information pertaining to the sentiment or opinion expressed in texts such as words, phrases, and documents. These methods draw from the aforementioned research areas. Notably, the assessment of sentiment polarity in a text is considered highly significant and is typically evaluated alongside its magnitude. Mejova (2009) also elucidates that sentiment detection can be viewed as a classification task, distinguishing texts as either objective or subjective, while polarity classification aims to gauge the intensity of polarity based on a chosen scale (e.g., positive, neutral, negative). Additionally, the discovery of the target of the opinion represents another important task.

In addition, Medhat et al. (2014) provide an overview of various sentiment classification methods. Firstly, they describe probabilistic classifiers such as the Naive Bayes Classifier, Bayesian Network, and Maximum Entropy Classifier. Next, they discuss linear classifiers including Support Vector Machines Classifiers and Neural Networks. Finally, decision trees and rule-based classifiers are presented. All of these methods fall under the category of machine learning approaches. The authors explain that in text classification, semi-supervised and unsupervised learning techniques are commonly employed due to the challenges associated with creating labeled training documents. Medhat et al. (2014) also highlight the use of lexicon-based approaches to determine the semantic orientation of text. These approaches rely on prepared sets of words that correspond to specific sentiments. Among the available options, dictionary-based and corpus-based approaches are mentioned.

Regarding the analysis of sentiment in tweets, Kharde and Sonawane (2016) conducted an evaluation of multiple techniques using a Twitter dataset. Their findings indicate that Naive Bayes and Support Vector Machine emerged as the most accurate methods. Consequently, the authors propose that these methods should serve as the baseline for machine learning approaches in sentiment analysis. Furthermore, they suggest that combining machine learning with lexicon-based methodologies would lead to superior results.

Qi and Shabrina (2023) arrived at comparable conclusions in their examination of Covid-19-related tweets sourced from individuals residing in England. Their study emphasizes the rapid evolution of language used in social media, thereby emphasizing the need for due consideration when employing lexicon-based methods. Additionally, the authors recommend utilizing sentiment classification outcomes in a regression model to explore the relationship between individuals' opinions and other pertinent variables.

The paper titled "Sentiment analysis in Twitter" by Martínez-Cámara et al. (2012) delves into the primary domains of investigation concerning Twitter, including polarity classification, political opinions, and event prediction, utilizing various Natural Language Processing (NLP) techniques. The paper's conclusion centers around the recognition of existing limitations in studying Twitter data, specifically highlighting three main obstacles: data sparsity, the presence of multiple languages, and the need for a well-defined scope when applying sentiment analysis.

Korkmaz et al. (2023) conducted a sentiment analysis to investigate the general perception of ChatGPT among Twitter users, two months after the chatbot's public release. The overall findings revealed that a significant majority of users (72%) expressed a positive orientation to-

wards ChatGPT, with only 6% adopting a neutral stance. On the other hand, negative emotions accounted for 22% of the sentiment, primarily encompassing "anger" and "fear." Furthermore, the intensity of positive emotions outweighed that of negative emotions. These results align with those obtained by Haque et al. (2022), indicating a high level of trust, excitement, and limited concern associated with the technology.

Regarding topic modeling, Kherwa and Bansal (2018) define it as a statistical technique employed to uncover the underlying semantic structure within extensive document collections. They identify commonly used methods such as Latent Semantic Analysis, Non-Negative Matrix Factorization, Probabilistic Latent Semantic Analysis, and Latent Dirichlet Allocation. The authors additionally highlight key applications of topic modeling, including scientific research, social network analysis, and software engineering.

Hong and Davison (2010) explored various topic modeling approaches using Twitter data and emphasized that the short length of tweets can present challenges in analysis. They suggest that aggregating short documents can partially overcome this limitation. Moreover, they demonstrate the usefulness of topic modeling in classification tasks.

Sanandres et al. (2020) exemplify the application of topic modeling on tweets to study and analyze conversations related to challenging social events in Colombia, while Sokolova et al. (2016) conducted a similar study focusing on Kenya. Both investigations aimed to gain deeper insights into the public perception of specific topics.

## 3 Methodology

### 3.1 Dataset description

The analysis conducted in this project relies on the dataset titled "#ChatGPT 1000 Daily Tweets," (Domingo, 2023) freely available online at Kaggle.com. This dataset is a collection of 66,375 tweets, each extracted based on the presence of at least one of the following words: ChatGPT, GPT3, or GPT4. The tweets have been compiled since the 3rd of April 2023 and are updated daily. Our analysis uses Version 39 of the dataset, published on the 14th of May 2023. The dataset comprises 20 distinct columns, each offering information related to the individual tweet. The full description for the dataset can be found in Appendix A. For this project, however, we have only utilized the 'text' column, representing the content of each tweet, the 'lang' column, indicating

the language in which the text is written, the ‘tweet\_created’ column showing the date the tweet was created, and the unique ‘tweet\_id’. The majority of the tweets are in a foreign language, with 20,545 English and 45,830 foreign tweets.

## 3.2 Pre-processing

### 3.2.1 Translation

The two main parts of the pre-processing of the data are the translation of foreign tweets into English, and the pre-processing of the tweets themselves. Before the translation of the foreign tweets, we have to make sure that all the data can be translated. As a number of tweet texts only include a URL, the first step was to remove such rows, and to remove URLs from tweets which contain both text and a link.

For the translation, Python’s *deep\_translator* package was used. The package is capable of translating to and from 133 languages using Google Translate’s API. We used the data’s ‘lang’ column to identify the languages which were used for the tweets. There were a few languages which were not supported by *deep\_translator*. Most of these were errors in the dataset (such as the language ‘1480.0’). As such, all of the tweets which had an unrecognized language were omitted. The only exception was the language ‘zh’. This language code is for the Chinese language, and was not recognized, as it appears in *deep\_translator* as ‘zh-CN’ (Simplified Chinese). For tweets in Chinese, the ‘lang’ column value was therefore changed to ‘zh-CN’.

After these steps, a function looked at whether the tweets were in English and translated them if they were not. The created values were outputted into a list, and the list was added to our dataframe as the column ‘translated\_text’.

### 3.2.2 Data Cleaning Process

The project involves two major processes - cleaning individual tweets with a function named *clean\_tweet* and implementing this function across the entire dataset.

The *clean\_tweet* function standardizes the text by converting it to lowercase, transforming emojis into words with similar meaning, removing URLs, user references, hashtags, and punctuation marks. It also filters out stopwords, and lemmatizes the words, reducing them to their root form. The outcome is a cleaner and standardized string of text, ready for analysis.

The *clean\_tweet* function is then systematically applied to the whole dataframe. This is

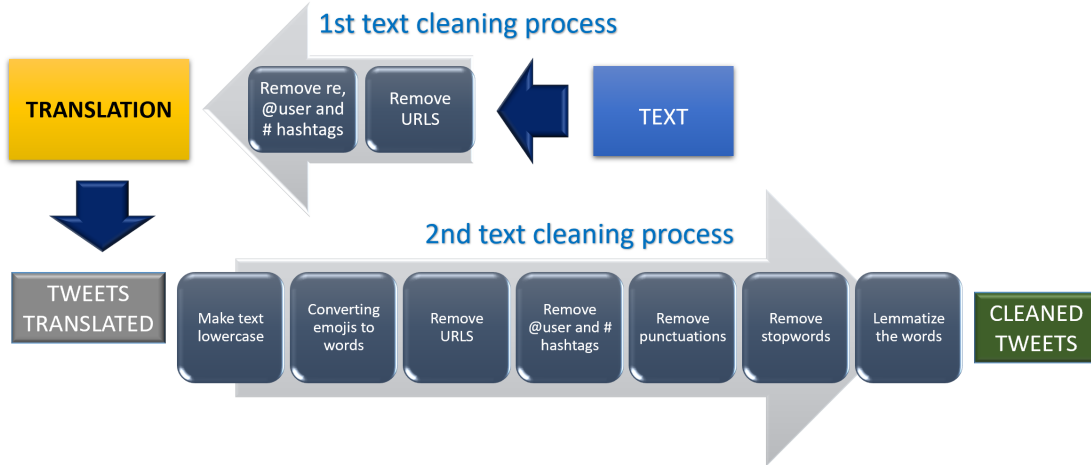


Figure 1: Text cleaning process.

achieved through a loop that processes each tweet, storing cleaned tweets in a list and keeping track of any errors that arise. Erroneous rows are subsequently dropped from the dataframe, ensuring a clean, preprocessed dataset.

After this comprehensive preprocessing stage, the cleaned tweets are added as a new column in the dataframe. Due to the way the pre-processing function works, a number of words without English meaning were included in the cleaned text. Most of these are names of foreign people. Because of this, the final step was to remove the words from the clean text which are not in the English dictionary. This cleaned dataset, free from noise and irregularities, is now ready for further analysis. Due to a number of preprocessing steps removing rows, our clean dataset has 39,885 rows of data. A visual representation of these processes can be seen in Figure 1.

### 3.3 Topic Modeling

The topic modeling process involves further data cleaning and feeding the data into a top2vec model. The main objective for this task is to find out the eight key underlying topics among the tweets. The top2vec model is selected among other probabilistic models, ie. LDA because of its efficiency in finding topics in a big row of documents that are specifically small texts like tweets.

As a preprocessing step, common words like “chatgpt”, “ai”, “openai” and “chatbot” have been removed from the cleaned tweets to achieve distinction and reduce noises from the topic words. After that the cleaned texts are transformed into a list of strings and fed into a top2vec model.

## Top2Vec Parameters

The list of strings is tuned with the parameters ‘min\_count = 0’ to include rare words, the speed has been set to ‘deep-learn’ to understand the corpus better. Also a pre-trained embedding ‘universal-sentence-encoder’ has been used to embed the text as other embeddings like “doc2vec” and “universal-sentence-encoder-multilingual” have demonstrated poor performance analyzing the topics words and producing a lot of noise.

## Umap Parameters

The “n\_neighbors” parameter in the umap parameters has been set to 50. It has been gradually increased from 5 to 50 to capture a more global structure as the tweets can have a lot of topics and it’s more in our interests to find out the overall theme of them. The number of components n\_components has been reduced to 2 for dimensionality reduction. The “metric” has been set to “cosine” as it produced the best results after trying out different metrics like “euclidean” and “manhattan”.

## HDBSCAN Parameters

The “min\_cluster\_size” has been gradually increased to 20 and “min\_samples” has been set to 50 after many trials to merge the smaller clusters into larger ones and control a balanced granularity. For the ‘cluster\_selection\_method’ the ‘eom’ method has been selected after comparing its performance with ‘leaf’ methods.

After running the model in the specified parameters, we received overall 131 topics. We have further reduced these topics to 10 topics using hierarchical topic reduction. These 10 topics will be further explored in the Discussion section.

## 3.4 Sentiment Analysis

### Polarity and Subjectivity

In the preliminary stages of our data exploration, we leveraged TextBlob, a Python library renowned for its sentiment analysis capabilities. The ‘sentiment’ attribute of TextBlob produces



a named tuple embodying sentiment in the form of polarity and subjectivity, applicable to the texts we had previously cleaned.

The TextBlob polarized scores are delivered as floating-point numbers within the range of  $[-1.0, 1.0]$ . The package operates by tokenizing the textual input and referencing a pre-existing internal dictionary. This dictionary assigns each word a polarity value, correlating to its perceived positivity or negativity. The overall polarity score for the text derives from the sum of these individual scores, divided by the sum of their absolute values.

Contrastingly, the subjectivity score, also a floating-point number, lies within  $[0.0, 1.0]$ . A score of 0.0 implies absolute objectivity, while 1.0 indicates high subjectivity. Hence, a score closer to 1.0 implies a subjective text, expressing personal opinions, emotions, or judgments. A score nearer to 0.0, however, suggests an objective text, primarily presenting factual information.

## **WordCloud**

The word clouds are created using the WordCloud package. The frequencies of aspects and sentiments are calculated using both the spaCy package and the Counter class from the collections module. The spaCy package is utilized to detect adjectives and nouns in our tweets database, while the Counter class helps us count the occurrences of these aspects and sentiments. The frequency of appearance of each aspect and sentiment is then used to determine the size and prominence of the words in the word cloud visualization. By visually representing the most frequently occurring adjectives and nouns, we gain insights into the prevalent aspects and sentiments expressed in our tweets data.

## **Word2Vec**

This project employs a two-stage process to conduct sentiment analysis and visualization on a cleaned tweet set, using tools like gensim, scikit-learn, and TextBlob. Initially, the Word2Vec model from gensim transforms each tweet into a corresponding vector, representing each unique word in a predefined vector space. Parameters like 'sentences', 'vector\_size', 'window', 'min\_count', 'workers', and 'max\_final\_vocab' help customize the model training.

After training, the word vectors are extracted, and their dimensionality is reduced using PCA, reducing the number of dimensions from 100 to 10. Further dimension reduction to 2 dimensions

is achieved using TSNE for better visualization. Lastly, sentiment polarity scores are calculated for each word using TextBlob, providing a numerical representation of sentiment ranging from -1.0 to 1.0, ready for further analysis or visualization.

## 4 Results and Discussion

### 4.1 Topic Modeling

Our initial 131 proposed topics were reduced to 10 using hierarchical topic reduction. These 10 topics were then further investigated to find their meanings.

- The first topic generated words like “Procrastination”, “Bot”, “Booming”, “Surpassing”, “Unstoppable”, “Obsolete”, “Peaked”, “Laziness”, “Earning”, etc. which are about the **Consequences** after the release of ChatGPT.
- The second topic generated words like Intelligent, “Futuristic”, “Sentience”, “Robot”, “Driverless”, “Humanoid” which are seemingly about the **Innovations** that are possible and were brought in after the introduction of ChatGPT.
- The third topic generated words like “Writing”, “Plagiarize”, “Proofreader”, “Transcribing”, “Writer”, “Essay”, “Poem”, etc which direct to be about using ChatBOT for various **Writing**.
- The fourth topic generated words like “What not”, “Bother”, “Overshadow”, “Amusingly”, “Jerk”, “Astonishingly”, that suggest tweets regarding **User Experiences** with ChatGPT.
- The fifth topic generated words like “Unanswered”, “Unhelpful”, “Poorly”, “Confused”, “Gibberish”, “Misinformation” which points towards the **Limitations & Fails** of ChatGPT.
- The sixth topic generated words like “Unethical”, “Professor”, “Prof”, “Student”, “Confidentiality”, “Misuse”, “Inappropriate”, etc which are mostly about **Academic & Ethical Concerns** regarding ChatGPT.
- The seventh topic generated words like “Prompt”, “Introductory”, “Beginner”, “Tutorial”, “Lifesaver”, “Nifty”, “Tip”, etc which suggests it is about the **Tips and Tricks** for beginners and users of ChatGPT.

- The eighth topic generated words like “Gratis”, “Premium”, “Free”, “Coupon”, ”Subscription”, ”Monetization”, etc which are about various **Subscription Services**.
- The ninth and tenth topic are generated words like “Wrench”, “Toolbox”, “Hammer”, “Cutter”, “Knife” (**Tool emojis**) and “Smile”, “Smiling”, “Smirking”, “Grinning”, “Frowning” (**Feeling emojis**) which are mostly transcriptions from the emojis which were transcribed for the sentiment analysis. Another approach could have been removing all the emojis before doing topic modeling for finding out more distinct topics.

## 4.2 Sentiment Analysis

According to our initial exploration with TextBlob polarity analysis, from all the tweets in our database since the launch of ChatGPT, 45% of users exhibited a positive attitude towards the implementation of ChatGPT. Similarly, 41% demonstrated a neutral sentiment about it, and only 14% showed a negative attitude toward ChatGPT (Figure 2).

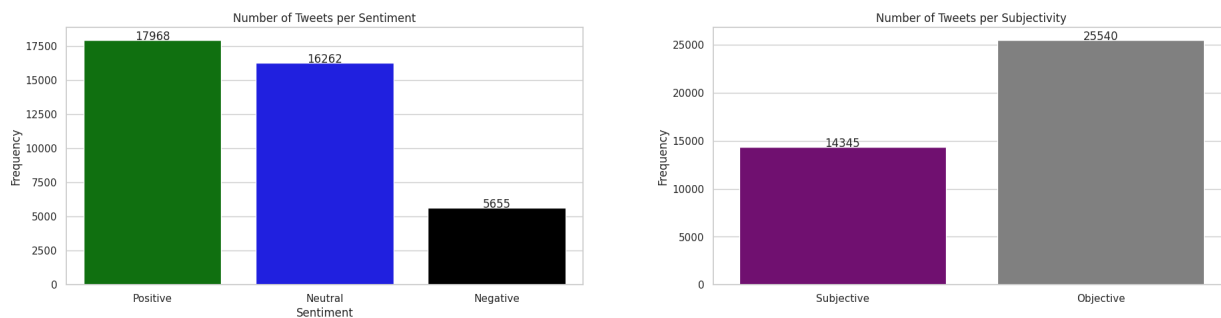


Figure 2: Number of Tweets per Sentiments and Subjectivity.

Also, using the same analysis but with the subjectivity parameter, we found that the conversation in tweets is only 36% subjective, based on opinions, judgments, and more. The rest, 64%, consists of objective texts. Therefore, we could hypothesize that the information shared on the social network is predominantly factual, detailing what ChatGPT is and what it does, with less personal sentiment or opinion expressed.

To contrast the previous findings, we conducted sentiment mapping using spaCy to detect adjectives and identify the most prominent ones in the tweets. Through this analysis, we determined the top 10 adjectives present in the tweet texts: new, free, artificial, prompt, many, good, first, last, and human (Appendix B.1). These adjectives were found to be highly representative of the

sentiments expressed in the tweets. Notably, the majority of these adjectives conveyed positive sentiments or highlighted factual information about ChatGPT.

By analyzing these sentiments and aspects, we gained valuable insights into the prevailing discussions surrounding ChatGPT. Our findings suggest that people generally hold a positive perception of the tool and often provide factual information about it. For a visual representation, refer to Appendix B.3.

We used Word2Vec to calculate the polarity of the word embeddings in our dataset, capturing their semantic meanings. Our findings indicate that, in general, words skew more towards the positive quadrants than the negative ones. This substantiates the overarching positive sentiment regarding the launch and use of ChatGPT. Figure 3, which highlights in the Quadrant I, displays words with a polarity above 0.5. These include "shocking", "good", "kind", "happiness", "intelligent" "impressive" but also "cruel", "horrible", "stupid" and "bad" to name a few. For a full view, the entire plot can be accessed in the Appendix B.4 section.

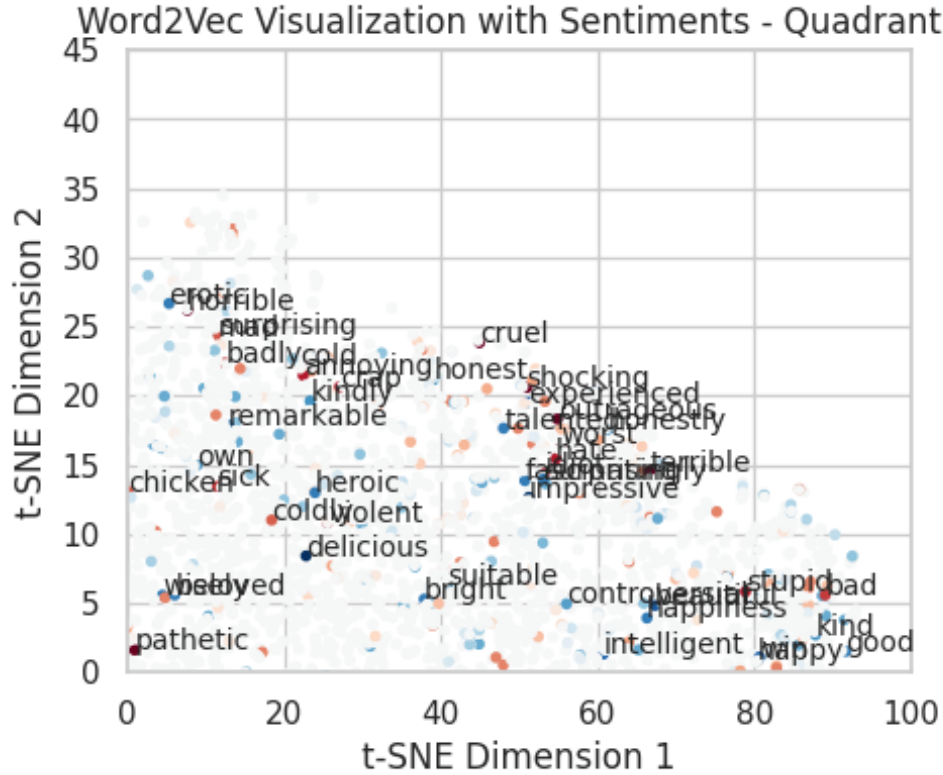


Figure 3: Sentiments with polarity higher than 0.5 in Quadrant I.

### 4.3 Sentiment Analysis over time

Our dataset contained around 1000 tweets from each day. We were interested to see if there is any significant trend during the examined period. For this, we looked at how the ratio of positive tweets compared to the number of positive and negative tweets changes. We excluded neutral sentiments, as we were interested about the changes in extreme feelings. Our findings for the daily and weekly sentiment ratios can be seen in Appendix B.5.

Overall, it can be observed a quite static trend over time. Weekly, the ratio of positive tweets lies between 0.75 and 0.80 (both values excluded) with week 16 and 19 showing the lowest result. But the variation is very limited and nothing suggests that it is related to any specific event or reason.

The daily sentiment analysis shows a deeper insight of people’s attitude towards ChatGPT, with slightly more fluctuation in terms of positive tweets ratio, having values going from 0.6 to over 0.85. Despite the lowest ratio of 0.60 nothing relevant happens and it can be concluded that people’s perception of ChatGPT remains positive. Nevertheless, it must be said that in just a month and a half wide changes in sentiment analysis would have been surprising, especially several months after the public release of the chatbot. Therefore it seems like that people have formed a fairly solid positive opinion of ChatGPT.

### 4.4 Sentiment Analysis of the Topics

From Figure 4, one can perceive that the overall impression for all the identified topics is positive, while the most negative sentiments are found in the Consequences section. Additionally, within the sections on Writing, User Experiences, and Limitations and Failures, negative sentiments are more prevalent than in other topics. Posts featuring face emojis are mostly positive, whereas posts with tool emojis carry more negative sentiments in comparison.

Generally some concern is shown but the general perception is overall positive, even among the topics which appear to be more critical such as those about limitations, consequences, and ethics. User Experiences shows some negative sentiments probably due to the very high number of users which slowed down the chatbot, while the Academic & Ethical Concerns negative tweets are likely related to plagiarism or similar issues. On the other hand, Innovation being the topic with the most positive positive ratio suggests trust in future applications of ChatGPT.

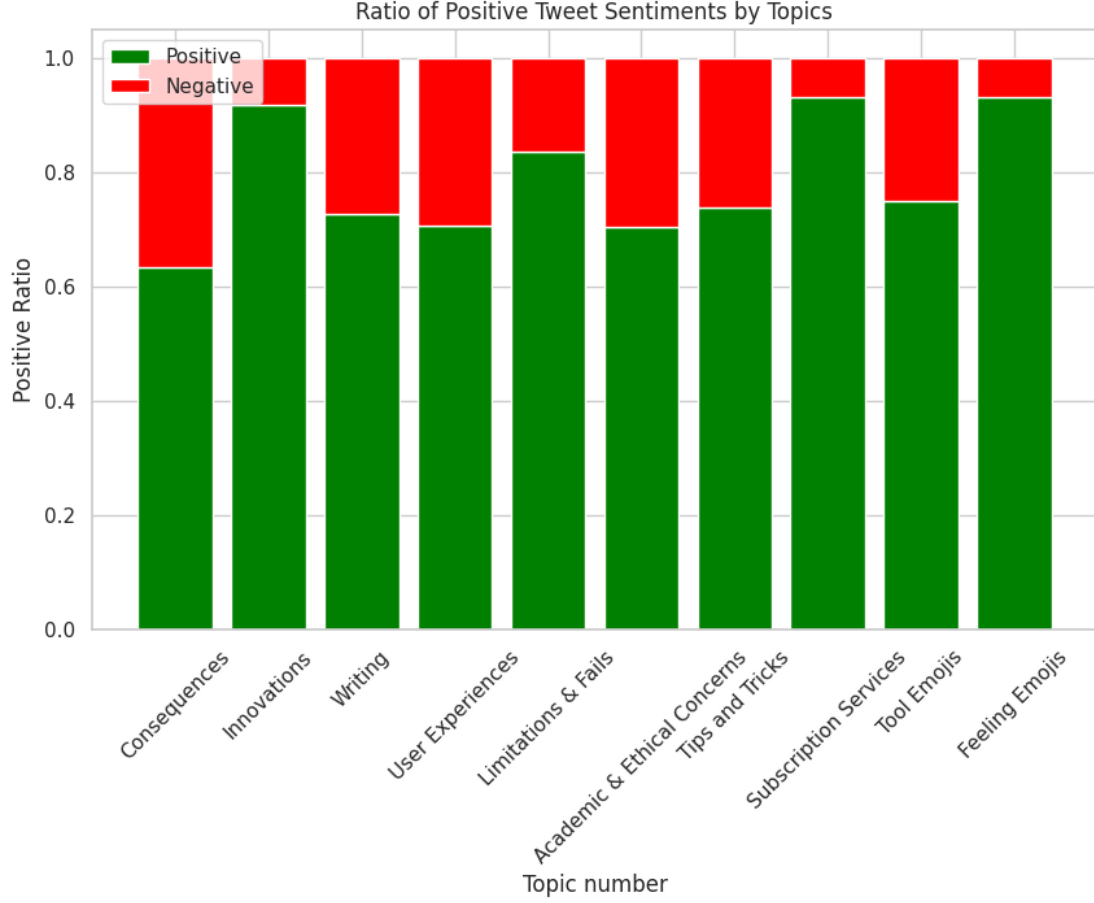


Figure 4: Sentiment scores for each topic.

## 5 Conclusion & Future Work

The project shows an overall good predisposition of ChatGPT users towards this new AI tool, with the majority of people perceiving it positively, while only a small minority appears to be concerned. These results strongly confirm previous findings obtained on two different datasets by Haque et al. (2022) and Korkmaz et al. (2023).

After conducting an exploration that involved implementing polarity and subjectivity measures, extracting adjectives, combining word embeddings with subjectivity measures, and utilizing top2vec for topic modeling, we can conclude that users generally exhibit a positive attitude towards the tool. The ratio of positive tweets to positive and negative tweets was 76%, meaning that for every negative sentiment, there are 3 positive ones.

The analysis was conducted including the emojis in tweets, after converting them into the words they represent. The choice is based on the idea that emojis are a relevant component in

social media conversation. However, despite many of them manifest emotions and are in the topic “Feeling emojis” providing additional elements for the project, the ones contained in “Tool emojis” may become misleading since, if considered out of context, their meaning could result unclear. Hence, implementing an analysis which excludes emojis could be an alternative possibility worth exploring, even though part of the content would be lost.

Due to the presence of only unlabeled data, the research project could only rely on lexicon-based methodologies. In fact, without any training, testing and validation data machine learning approaches cannot be implemented. Therefore, in the future it could be useful to explore the possibilities offered by such methods, such as the Naive Bayes and Support Vector Machine (Kharde and Sonawane, 2016).

Additionally, it should be kept in mind that the used dataset covers a limited time frame (just one month and a half) and it considers a period of time quite subsequent to the release of ChatGPT. Hence, further research may investigate people’s attitude towards this chatbot since it became openly available in order to see also how sentiment has been changing over time.

The dataset also contained information about users’ locations. However, a large portion of this data was missing or erroneous, therefore it was not examined in this paper. Future researchers may look at the actual location data in order to examine continent or country-wide differences.

# References

- Domingo, E. (2023). Chatgpt 1000 daily tweets.
- Haque, M. U., Dharmadasa, I., Sworna, Z. T., Rajapakse, R. N., and Ahmad, H. (2022). "i think this is the most disruptive technology": Exploring sentiments of chatgpt early adopters using twitter data.
- Hong, L. and Davison, B. D. (2010). Empirical study of topic modeling in twitter. In *Proceedings of the First Workshop on Social Media Analytics*. ACM.
- Kharde, V. A. and Sonawane, S. (2016). Sentiment analysis of twitter data: A survey of techniques. *International Journal of Computer Applications*, 139(11):5–15.
- Kherwa, P. and Bansal, P. (2018). Topic modeling: A comprehensive review. *ICST Transactions on Scalable Information Systems*, 0(0):159623.
- Korkmaz, A., Aktürk, C., and Talan, T. (2023). Analyzing the user's sentiments of ChatGPT using twitter data. *Iraqi Journal for Computer Science and Mathematics*, pages 202–214.
- Martínez-Cámara, E., Martín-Valdivia, M. T., Ureña-López, L. A., and Montejo-Ráez, A. (2012). Sentiment analysis in twitter. *Natural Language Engineering*, 20(1):1–28.
- Medhat, W., Hassan, A., and Korashy, H. (2014). Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal*, 5(4):1093–1113.
- Mejova, Y. (2009). Sentiment analysis: An overview. *University of Iowa, Computer Science Department*.
- Qi, Y. and Shabrina, Z. (2023). Sentiment analysis using twitter data: a comparative application of lexicon- and machine-learning-based approach. *Social Network Analysis and Mining*, 13(1).
- Sanandres, E., Abello, R., and Madariaga, C. (2020). Topic modeling of twitter conversations: The case of the national university of colombia. In *Studies in Classification, Data Analysis, and Knowledge Organization*, pages 241–251. Springer International Publishing.
- Sokolova, M., Huang, K., Matwin, S., Ramisch, J., Sazonova, V., Black, R., Orwa, C., Ochieng, S., and Sambuli, N. (2016). Topic modelling and event identification from twitter textual data.



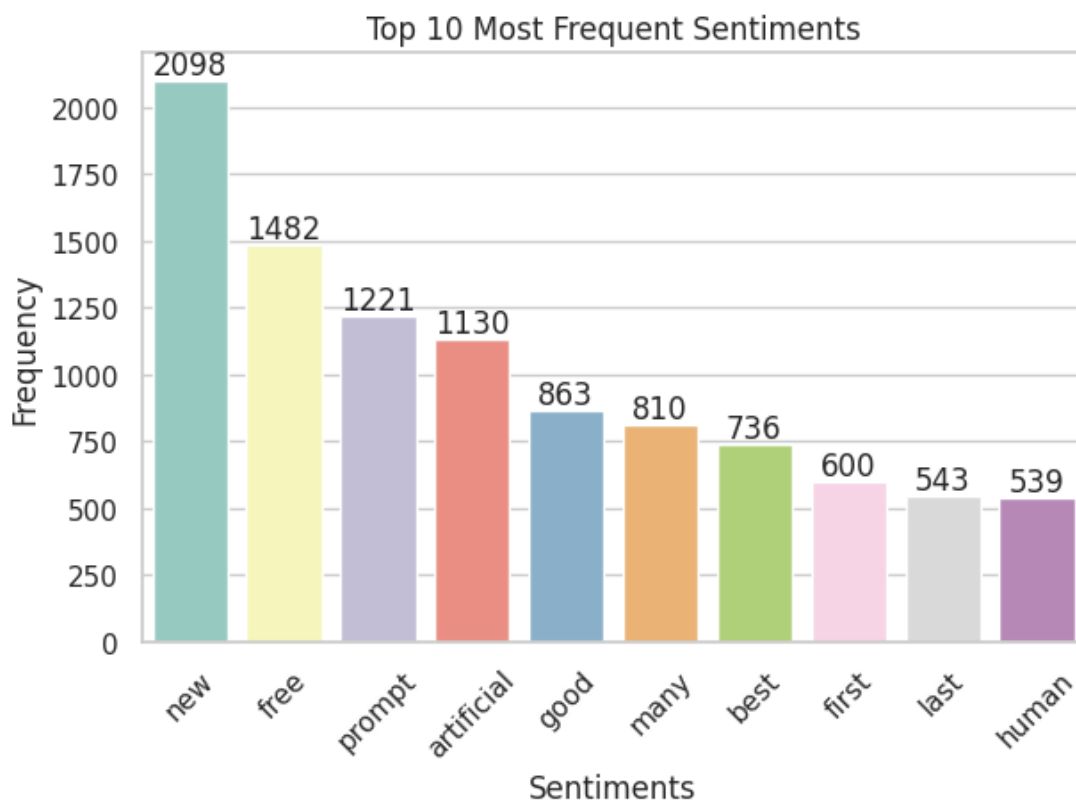
## A Data overview

Column	Description
tweet_id	key field
tweet_created	date of creation
tweet_extracted	date of extraction of the tweet by the creator of the dataset
text	corpus of each tweet
lang	language in which the text written
user_id	Unique identifier for each Twitter user who posted a tweet
user_name	Display name of the Twitter user who posted the tweet
user_username	Username (handle) of the Twitter user who posted the tweet
user_location	The location as provided by the user in their profile, usually the user's country
user_description	A brief personal description written by the user in their Twitter profile
user_created	Date and time when the Twitter user's account was created
user_followers_count	The number of followers that the user has at the time the tweet was extracted
user_tweet_count	The total number of tweets (including retweets) posted by the user at the time the tweet was extracted
user_verified	Boolean value indicating whether the user is verified by Twitter (True) or not (False)
source	The device or application used by the user to post the tweet
retweet_count	The number of times the tweet has been retweeted at the time of extraction
like_count	The number of likes the tweet has received at the time of extraction
reply_count	The number of replies the tweet has received at the time of extraction
impression_count	The number of times the tweet was viewed (impressions) at the time of extraction

## B Sentiment Analysis

### B.1 Most Frequent Sentiments

We employed spaCy for sentiment analysis to detect adjectives and identify the most prominent ones in the tweets. Through this analysis, we identified the top 10 adjectives prevalent in the tweet texts.



## B.2 200 most prominent nouns

Utilizing the same method as before, we detected the most prominent nouns in the tweets. The following WordCloud showcases the top 200 of these.

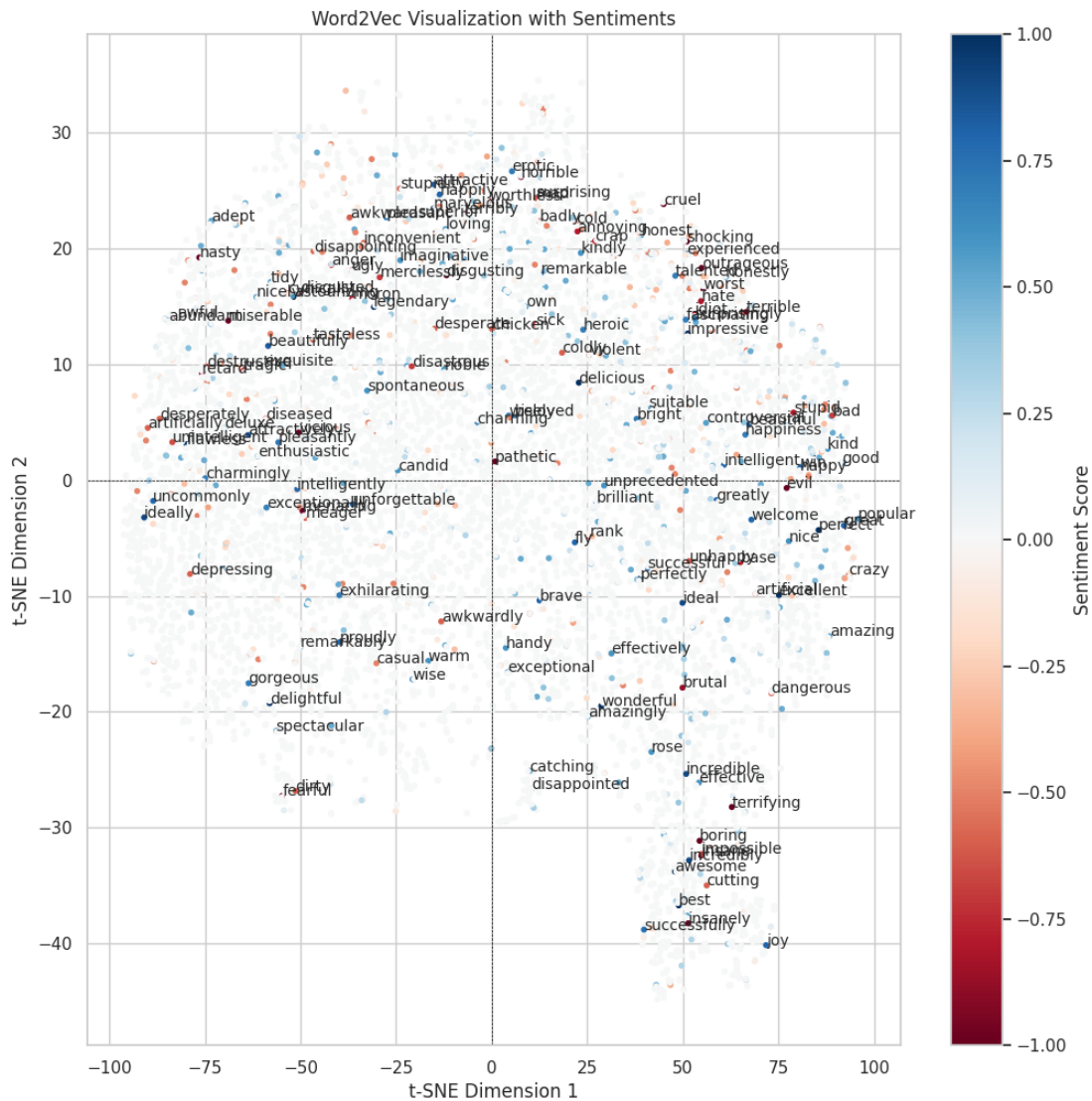


### B.3 Most prominent sentiments



## B.4 Word embeddings with high sentiment score visualization

In the following scatter plot, we've used Word2Vec to depict the most prominent words, combined with TextBlob polarity measures. The color represents the sentiment of the word: blue denotes negative sentiment (-1), white indicates neutral sentiment (0), and red symbolizes positive sentiment (1), considering we employed the 'RdBu' (Red-Blue) colormap. The more intense the color, the stronger the sentiment of the word. Words situated close together on the plot are used in similar contexts within the tweets. Conversely, words that are placed far apart are used in dissimilar contexts and are likely to carry different meanings or usages.



## B.5 Sentiment Analysis over time

