

# Overview of Machine Learning Algorithms for Detecting Microaggression in Written Text

Asif Tareque, Harshith Hullakere Siddegowda, Denster Joseph Frank, Nicole Lee, Rezza Moieni

Diversity Atlas, Melbourne, Australia

Email: nicole.lee@diversityatlas.io, rezza.moieni@diversityatlas.io

**How to cite this paper:** Tareque, A., Siddegowda, H. H., Frank, D. J., Lee, N., & Moieni, R. (2024). Overview of Machine Learning Algorithms for Detecting Microaggression in Written Text. *Open Journal of Social Sciences*, 12, 347-358. <https://doi.org/10.4236/jss.2024.127025>

**Received:** December 1, 2023

**Accepted:** July 19, 2024

**Published:** July 22, 2024

Copyright © 2024 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## Abstract

Microaggressions are brief, daily verbal, behavioral, or environmental actions that convey negative, demeaning, or hostile racial undertones. These can be unintentional and often go unnoticed by the offender. They can have significant impacts on the mental health of the victims, leading to stress, low self-esteem, and feelings of invalidation. This research aims to detect microaggressions in written communication using machine learning. The study tackles the problem of data scarcity and lacking annotated data on microaggression, by collecting text data from microaggressions.com, ChatGPT, Reddit and office workplaces, and annotating the data using GPT 3.5 language model. Multiple machine learning algorithms were used to detect microaggressive language in text and were evaluated across proper metrics. Long Short-Term Memory (LSTM) with BERT embeddings was found to be the most stable model in detecting microaggression. It advances the field of microaggression detection with the leveraging of deep learning techniques, which could be potentially expanded to eliminate microaggression in texts.

## Keywords

Natural Language Processing, Inclusivity, Microaggression, Diversity, Machine Learning, Deep Learning

## 1. Introduction

Microaggressions are short, everyday verbal, behavioral, or environmental actions, whether deliberate or not, that convey negative, demeaning, or hostile racial undertones (Sue et al., 2007). Initially, it was believed that victims existed exclusively in racial minorities or people of colour, however, microaggression can happen to anyone based on race, gender, sexual orientation or any other

protected characteristic (Sue, 2010).

The subtle nature of microaggression makes them insidiously dangerous as offenders often commit them unintentionally and are unaware of the consequences for the victim. Microaggressions have been compared to deadly carbon monoxide gas which is potentially lethal but undetectable (Sue & Sue, 2003) when offenders are not conscious of such microaggressive behavior they may be less inclined to correct their behaviour and may even justify it to themselves. Research, however, suggests microaggressions impact the victims greatly and thus are worthy of further research.

Correlation and regression analysis reveal the link between racial microaggressions and mental health and concludes that the more microaggressions a person experiences, there are associated poor mental health outcomes. Notably, microaggressions were significantly correlated with negative mental health outcomes, particularly in depression, lack of positive mood, and lack of behavioral control (Nadal et al., 2014). Furthermore, while individual microaggressions might seem minor in the wider scheme of things, their cumulative effect is potentially substantial. Over time such effects may lead to stress, lack of self-esteem and loneliness. Furthermore, microaggressions may cause victims a sense of invalidation within the world, thus further impacting mental health outcomes. Lastly, within the context of clinical therapy, the occurrence of microaggressions can deteriorate the therapeutic alliance between client and therapist (Sue et al., 2007).

Utilizing the Racial and Ethnic Microaggressions Scales (REMS) and the Mental Health Inventory (MHI) it was found that racial microaggressions were experienced significantly more by people with general mental health issues. Furthermore, it was found that factors such as geographic location, education, and age can also be factors of experiencing microaggressions (Nadal, 2018). Microaggressions can also be termed as algorithmic bias in language models, as a result, companies hold an ethical responsibility for the inflicted damage caused by their algorithms (McClure & Wald, 2022).

While a lot of work has been done in classifying toxic written text, such as comments on social media (Aken et al., 2018), little research has been completed on using machine learning to detect microaggressions. Therefore, this study aims to answer the following research questions:

- 1) How to detect microaggression in written communication using Machine learning?
- 2) Can synthetic data tackle the lack of microaggression data?

## 2. Literature Review

Natural Language Processing (NLP) has been used previously to overcome some of the challenges that the diversity and inclusion field faces such as promoting gender equity (Raichur, Lee, & Moieni, 2023). Previous research has also found that there is a lack of work completed about detecting microaggressions and

proposes an automated racial microaggression detection tool using Random Forest and IBk (KNN classifier of Weka Software) classification algorithms (Ali et al., 2020) where the results seem consistent in detecting non-microaggression and promising in terms of detecting racial microaggression. However, it only performs well when explicitly microaggressive text is the input. It does not perform well with general texts such as newspaper articles, however. The shortcoming arises from the lack of data and class imbalance, as the data set had almost two-thirds of the text labelled microaggressive.

Furthermore, studies have been done to research the effectiveness of Machine Learning models in detecting microaggressions in various contexts such as workplaces, social media, and general conversations. These studies show how such models could detect microaggressions in scripted TV shows when trained on real life conversation (Ngueajio et al., 2023). This research implements a Support Vector Machine (SVM) with N-grams for feature representation in one model and Robustly Optimized Bidirectional Encoder Representation from Transformers (ROBERT) for context-based feature representation for another model. The paper concludes that the contextual model simply outperforms the model that uses N-grams for feature representation and that models trained on real-life conversations were able to detect microaggression in scripted TV settings at equal rates.

Moreover, unsupervised machine learning algorithms have been used in working to detect microaggressions (Ògúnremí, Basile, & Caselli, 2022). The research shows that inherent biases present in pre-trained word embedding could be used to pinpoint subtle, offensive language patterns, particularly microaggression. While unsupervised algorithms do not require labeled data to operate, the algorithm is not able to detect implicit, othering phrasing commonly associated with microaggressions, such as “your kind” or “you people”. Additionally, the challenge of polysemy, where a word can have multiple meanings, presents another obstacle to unsupervised machine learning. Hence, we propose a supervised approach to address these challenges.

Nonetheless, this research highlights a key challenge in microaggression research: a lack of annotated data. The challenge lies in the lack of real-world data specifically annotated as microaggression (Breitfeller et al., 2019) as opposed to general aggression or hate speech. While the traditional solution involves crowd sourcing data through platforms such as MTurk, in the era of AI, we propose annotation and synthetic data generation using GPT to overcome this data shortage.

Recent research has shown GPT achieves a 70% accuracy rate for content moderation tweets data, 81% for news articles data and 83% for US Congress tweets data. In terms of GPT’s intercoder agreement, ChatGPT performs significantly better: MTurk averages 56%, trained annotators 79%, ChatGPT (temperature = 1) 91%, and ChatGPT (temperature = 0.2) achieves a remarkable 97%. ChatGPT’s temperature parameter controls the degree of randomness of the output (Gilardi, Alizadeh, & Kubli, 2023). Thus, large generative language mod-

els could be considered in performing data annotation with significant reliability.

Furthermore, with the rise to End-to-end machine learning while working with NLP where raw data is inserted as input, as opposed to manually engineering features and training both encoder and model at once. One way of achieving this is to implement Bidirectional Encoder Representation (BERT) as an embedding layer on top of which another neural network architecture is set (Li et al., 2019). Word2Vec or GloVe-based embedding layers generate a context-independent representation for each token and the BERT embedding layer calculates the token-level representation using the information from the overall sentence provided as input. Therefore, using such an approach removes the overall challenge of feature engineering required for such a task. Another study (Miaschi & Dell'Orletta, 2020) compares the probing score of contextual BERT word embeddings and non-contextual word2vec. BERT captures features related to the basic text and sentence structure, while Word2vec is better at predicting word formation and sentence structure details.

### 3. Methodology

Figure 1 shows an outline of this paper's methodology which is elaborated in details in the upcoming subsections.

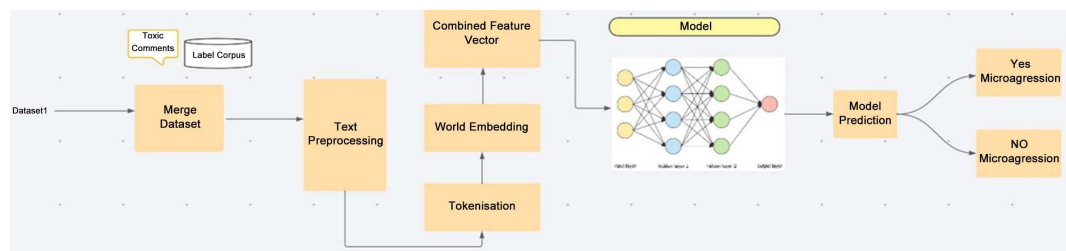


Figure 1. Overall architecture on data collection, data annotation, model evaluation, model selection, and microaggression detection.

#### 3.1. Data Collection

To compile a diverse and comprehensive dataset for the research, we employed various data collection techniques from different sources as listed below. By employing diverse methods, we aimed to construct a well-rounded dataset that encompasses a wide range of perspectives and situations, enabling a comprehensive analysis for our research objectives.

##### 1) Synthetic Data Generation from Chat GPT (GPT-16k Turbo):

Due to ethical considerations and limitations arising from the updated policies of Chat GPT, the study resorted to generating synthetic data using the GPT-16k Turbo model. This, however, gives rise to challenges pertaining to both ethical concerns and the quantity of generated data. To mitigate these issues, the generation process was conducted in batches. Ultimately, approximately 1700 data points were collected, some which are shown in **Appendix**.

##### 2) Reddit API Data Collection:

To incorporate real-world perspectives and opinions on controversial topics, we utilized the Reddit API. We focused on a selection of topics known for their contentious nature, including social justice, Black Lives Matter, feminism, LGBTQI+ issues, Asian American, Native American, Latino America, disability, Muslim Lounge, Jewish, Climate Action, Technology, Science, Mental Health, Personal Finance, Parenting, Travel, Books, Fitness, and Art. From each of these categories, we extracted the first 10 hot topics. In total, 725 data points were collected. Subsequently, the data underwent a thorough cleaning and pre-processing process.

### 3) *Data Scraping from microaggressions.com:*

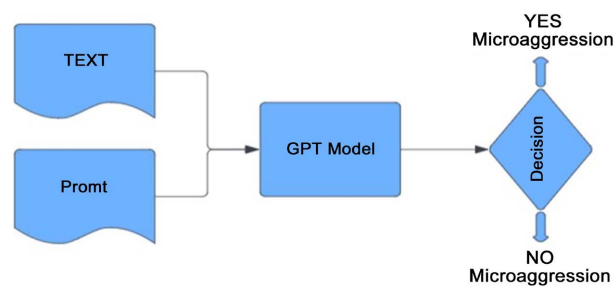
Employing web scraping techniques, we gathered data from microaggressions.com ([Microaggressions Project, 2023](#)). This source provided valuable insights into instances of microaggressions, which are often critical in understanding subtle forms of discrimination and bias.

### 4) *Email Conversations from Office Workplace.*

Additionally, the data was sourced from email conversations within office workplaces ([Civil Research Data, 2018](#)). This particular mode of communication provides a unique perspective on professional interactions and sheds light on various workplace dynamics.

## 3.2. Data Annotation Using GPT 3.5

Following the initial data collection phase, deliberation arose regarding the choice between employing human annotators or utilizing Language Learning Models (LLMs) for data annotation as shown in [Figure 2](#). Subsequent to an extensive review of the existing literature ([Gilardi, Alizadeh, & Kubli, 2023](#)), it was conclusively determined that Chat GPT exhibited superior performance compared to human annotators in the task of text annotation. Consequently, the GPT-3.5 Turbo model with a capacity of 16k tokens were employed for the annotation of data into binary classifications of “Yes Microaggression” and “No Microaggression”.



**Figure 2.** Architecture on data annotation using Language Learning Models (LLMs).

## 3.3. Data Preprocessing

In the data cleaning phase, a series of NLP techniques were applied to enhance the quality and consistency of the dataset. These procedures encompassed the

following steps:

1) Converting to Lowercase: All text entries were transformed into lowercase, ensuring uniformity and simplifying subsequent analyses.

2) Whitespace Removal: Extraneous spaces within the text were systematically eliminated, further streamlining the dataset.

3) Special Character Removal: Any non-alphanumeric characters were excised from the text, eliminating potential sources of noise or irregularities.

4) URL Elimination: Uniform Resource Locators (URLs) were systematically removed to prevent them from influencing subsequent analyses.

5) Emoji Removal: Emoticons and other non-textual symbols were purged from the dataset, focusing the analysis exclusively on textual content.

6) Non-English Language Exclusion: Texts not in English were identified and subsequently excluded from the dataset, ensuring linguistic homogeneity.

7) Lemmatization: Lemmatization was preferred over stemming due to its ability to produce linguistically valid root words. Unlike stemming, which often results in non-standard or even non-existent words, lemmatization preserves the semantic integrity of the text. This ensures that the derived root words maintain their meaningfulness within the context of the language.

8) Retention of Stop Words: As per the findings of our comprehensive literature review in section 2, the decision was made to retain stop words. Removal of these common linguistic elements can lead to a loss of crucial context and nuance, potentially impeding subsequent analyses.

### 3.4. Model Selection

This research experiments with four classification techniques, including both non-neural and deep neural network:

1) Logistic regression (LR):

LR is a traditional statistical tool has been gaining attention in the realm of machine learning, particularly in the domain of text classification (Zhang et al., 2003; Genkin, Lewis, & Madigan, 2007; Ifrim, Bakir, & Weikum, 2008).

2) Support Vector Machine (SVM):

SVM is a linear model that is commonly applied in binary classification (Steinwart & Christmann, 2008). It served as essential baseline model, providing valuable insights into the task's initial complexities and setting the foundation for subsequent investigations in this study.

3) Long Short-Term Memory (LSTM) with BERT embedding layer:

LSTM is a deep neural network that is applicable to text classification and prediction (Nowak, Taspinar, & Scherer, 2017). While the BERT model is a language model capable of capturing nuanced contextual information (Devlin et al., 2019). The combination of both is a sophisticated approach in NLP to understand the subtle nuances of language and context (Pandey & Singh, 2023), which is a critical requirement for accurately identifying microaggressions in text.

4) Gated Recurrent Units (GRU):

GRU is another deep neural network with the capability of learning long sequences of text. It is gaining popularity by its simplicity and fewer parameters compared to the LSTM models (Zulqarnain et al., 2020).

### 3.5. Evaluation Metrics

This study incorporates the assessment of accuracy and loss that is relevant in the context of deep learning models. The degree of fitting of each model was assessed in terms of training and testing accuracy. Overfitting exists when the model performs well on training data and performs badly on testing data. Underfitting, on the other hand, exists when the model performs badly on training data and performs well on testing data. A model is considered stable when it generalises well with both seen data on training data and unseen data on testing data, with no overfitting or underfitting (Ying, 2019).

In addition, criteria to evaluate the models involve three key metrics, which were assessed on both microaggressive (1) and non-microaggressive (0) text:

- 1) Precision: The measure of positive observations that are correctly predicted from the total predicted observations in the class.
- 2) Recall: The fraction of class from the total correctly classified once.
- 3) F1-score: The harmonic balance of precision and recall (Hossin & Sulaiman, 2015).

While comparing, the F1 score will suggest the better model and standardize comparison with relevant previous works. We expect a higher recall from the model in the microaggressive class and high precision in the non-microaggressive class. This preference arises from the importance of correctly detecting microaggressions and not mislabeling microaggressive texts, as failing to do so could result in catastrophic consequences, especially when numerous microaggressive texts go undetected.

Together, these metrics provide a multifaceted evaluation framework, enabling a thorough understanding of the model's performance over the classification task and model optimization.

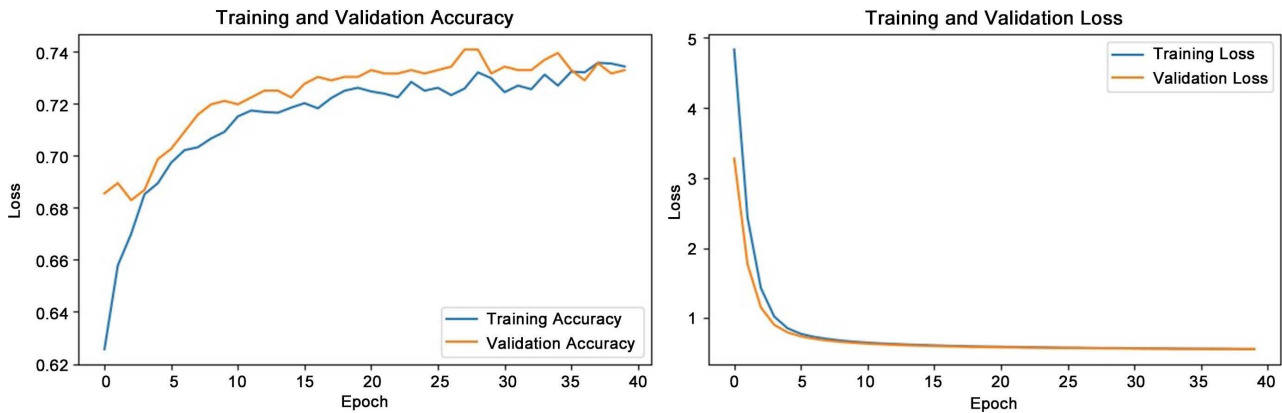
## 4. Results

All four models have roughly the same test accuracies, as described in **Table 1**, from which GRU under-performed at the lowest accuracy value of 72% on the test data. Meanwhile, LSTM has a training accuracy of 73% and test accuracy of 74%, showing no overfitting and is consistent with both seen and unseen data, which appears to be the most stable model.

**Table 1.** Training and test accuracy value of different models.

Model	Training accuracy	Test accuracy
LR	0.8313	0.7505
SVM	0.9852	0.7474
LSTM	0.7320	0.7408
GRU	0.7929	0.7197

Furthermore, the overall performance of LSTM over multiple epochs is presented in **Figure 3**. The figure on the left shows the expected behaviour wherein both training and validation accuracy exhibit steady improvement before eventually plateauing. Meanwhile, in the figure on the right, convergence loss occurs when the training and validation loss decreases steadily over epochs before eventually reaching a point where losses are not reduced any further. Therefore, we are assured that the model has the capability of learning as much information as possible from available data without overfitting.



**Figure 3.** Trend of training and validation accuracy (left) and training and validation loss (right) over multiple epochs.

Comparing the F1 scores of each model, all models perform better with higher F1 scores when classifying and detecting non microaggression. A summarized evaluation results could be found in **Table 2**. Focusing on the LSTM model, when it predicts the text as non-microaggressive, the precision value of “No Microaggression” indicates a 72% chance that the text is genuinely non-microaggressive. On the other hand, the recall value of the “Microaggression” class indicates the model could correctly detect 59% of the microaggressions. Thus, the LSTM model is more reliable in predicting text with no microaggression.

Nonetheless, the recall value for detecting microaggressions in this study is higher than the previous research of Ali (Ali et al., 2020). A key factor contributing to this improved performance could be the use of a less imbalanced dataset, achieved through the generation of synthetic data and data annotation using the language model GPT.

**Table 2.** Evaluation results of different models across three metrics.

Model	No microaggression			Yes Microaggression		
	Precision	Recall	F1 score	Precision	Recall	F1 score
LR	0.72	0.90	0.80	0.82	0.57	0.67
SVM	0.76	0.78	0.77	0.73	0.70	0.71
LSTM	0.72	0.84	0.77	0.74	0.59	0.66
GRU	0.72	0.84	0.78	0.75	0.59	0.66

## 5. Conclusion

This study hypothesized that machine learning and artificial intelligence have the potential to overcome barriers to diversity and inclusion by detecting microaggressions in texts. It overcomes the challenge of annotated data regarding microaggressions by using annotation through LLMs which resulted in improved performance due to better data balance. It reveals that LSTM model exhibited the best performance in microaggression detection.

While detection is one key aspect of this problem of microaggression, further study could be further enhanced by the use of language models to improve text by perhaps paraphrasing texts to be non-microaggressive where the existing models could be put as a feedback loop for the future generative model. Thus, we leave out this research question for future work to tackle microaggression, “If Artificial Intelligence (AI) can reduce the microaggression in written communication?”.

## Acknowledgements

This research was supported by Diversity Atlas and their provision of data has been instrumental in shaping the findings of this study. We would like to express our special thanks to the RMIT University for supporting the internship program that allowed the authors to conduct this research endeavor.

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

- Aken, B. V., Risch, J., Krestel, R., & Löser, A. (2018). *Challenges for Toxic Comment Classification: An In-Depth Error Analysis*. In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)* (pp. 33-42). Association for Computational Linguistics.
- Ali, O., Scheidt, N., Gegov, A., Haig, E., Adda, M., & Aziz, B. (2020). Automated Detection of Racial Microaggressions Using Machine Learning. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)* (pp. 2477-2484). IEEE. <https://doi.org/10.1109/SSCI47803.2020.9308569>
- Breitfeller, L., Ahn, E., Jurgens, D., & Tsvetkov, Y. (2019). Finding Microaggressions in the Wild: A Case for Locating Elusive Phenomena in Social Media Posts. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (pp. 1664-1674). Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-1176>
- Civil Research Data (2018). *Data.json*. <https://figshare.com/articles/dataset/data,json/7376747>
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Lin-*

- guistics: Human Language Technologies* (pp. 4171-4186). Association for Computational Linguistics. <https://doi.org/10.18653/v1/N19-1423>
- Genkin, A., Lewis, D. D., & Madigan, D. (2007). Large-Scale Bayesian Logistic Regression for Text Categorization. *Technometrics*, 49, 291-304. <https://doi.org/10.1198/004017007000000245>
- Gilardi, F., Alizadeh, M., & Kubli, M. (2023). ChatGPT Outperforms Crowd-Workers for Text-Annotation Tasks. *Proceedings of the National Academy of Sciences of the United States of America*, 120, e2305016120. <https://doi.org/10.1073/pnas.2305016120>
- Hossin, M., & Sulaiman, M. N. (2015). A Review on Evaluation Metrics for Data Classification Evaluations. *International Journal of Data Mining & Knowledge Management Process*, 5. <https://doi.org/10.5121/ijdkp.2015.5201>
- Ifrim, G., Bakir, G., & Weikum, G. (2008). Fast Logistic Regression for Text Categorization with Variable-Length N-Grams. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 354-362). Association for Computing Machinery. <https://doi.org/10.1145/1401890.1401936>
- Li, X., Bing, L., Zhang, W., & Lam, W. (2019). Exploiting BERT for End-to-End Aspect-Based Sentiment Analysis. In *Proceedings of the 5th Workshop on Noisy User-Generated Text (W-NUT 2019)* (pp. 34-41). Association for Computational Linguistics. <https://doi.org/10.18653/v1/D19-5505>
- McClure, E. & Wald, B. (2022). Algorithmic Microaggressions. *Feminist Philosophy Quarterly*, 8, Article 5. <https://doi.org/10.5206/fpq/2022.3/4.14276>
- Miaschi, A., & Dell'Orletta, F. (2020). Contextual and Non-Contextual Word Embeddings: An In-Depth Linguistic Investigation. In *Proceedings of the 5th Workshop on Representation Learning for NLP* (pp. 110-119). Association for Computational Linguistics. <https://doi.org/10.18653/v1/2020.repl4nlp-1.15>
- Microaggressions Project (2023). *Microaggressions in Everyday Life*. <https://www.microaggressions.com/>
- Nadal, K. L. (2018). *Microaggressions and Traumatic Stress: Theory, Research, and Clinical Treatment*. American Psychological Association. <https://doi.org/10.1037/0000073-000>
- Nadal, K. L., Griffin, K. E., Wong, Y., Hamit, S., & Rasmus, M. (2014). The Impact of Racial Microaggressions on Mental Health: Counseling Implications for Clients of Color. *Journal of Counseling & Development*, 92, 57-66. <https://doi.org/10.1002/j.1556-6676.2014.00130.x>
- Ngueajio, M. K., Hernandez, I., Cornett, K., Washington, G., & Parsons, D. (2023). Towards Identification of Microaggressions in Real-Life and Scripted Conversations, Using Context-Aware Machine Learning Techniques. <https://openreview.net/forum?id=z7FfWq2iaW4>
- Nowak, J., Taspinar, A., & Scherer, R. (2017). LSTM Recurrent Neural Networks for Short Text and Sentiment Classification. In L. Rutkowski, M. Korytkowski, R. Scherer, R. Tadeusiewicz, L. Zadeh, & J. Zurada (Eds.), *Artificial Intelligence and Soft Computing* (pp. 553-562). Springer International Publishing. [https://doi.org/10.1007/978-3-319-59060-8\\_50](https://doi.org/10.1007/978-3-319-59060-8_50)
- Ògúnrèmi, T., Basile, V., & Caselli, T. (2022). Leveraging Bias in Pre-Trained Word Embeddings for Unsupervised Microaggression Detection. *Italian Journal of Computational Linguistics*, 8. <https://doi.org/10.4000/ijcol.1066>
- Pandey, R., & Singh, J. P. (2023). BERT-LSTM Model for Sarcasm Detection in Code-Mixed Social Media Post. *Journal of Intelligent Information Systems*, 60, 235-254.

<https://doi.org/10.1007/s10844-022-00755-z>

Raichur, A., Lee, N., & Moieni, R. (2023). A Natural Language Processing Approach to Promote Gender Equality: Analysing the Progress of Gender-Inclusive Language on the Victorian Government Website. *Open Journal of Social Sciences*, 11, 513-529.

<https://doi.org/10.4236/jss.2023.119033>

Steinwart, I., & Christmann, A. (2008). *Support Vector Machines*. Springer Science & Business Media.

Sue, D. W. (2010). *Microaggressions in Everyday Life: Race, Gender, and Sexual Orientation*. John Wiley & Sons.

Sue, D. W., & Sue, D. (2003). *Counseling the Culturally Diverse: Theory and Practice*. John Wiley & Sons.

Sue, D. W., Capodilupo, C. M., Torino, G. C., Bucceri, J. M., Holder, A., Nadal, K. L., & Esquilin, M. (2007). Racial Microaggressions in Everyday Life: Implications for Clinical Practice. *American Psychologist*, 62, 271-286.

<https://doi.org/10.1037/0003-066X.62.4.271>

Ying, X. (2019). An Overview of Overfitting and Its Solutions. *Journal of Physics: Conference Series*, 1168, Article ID: 022022.

<https://doi.org/10.1088/1742-6596/1168/2/022022>

Zhang, J., Jin, R., Yang, Y., & Hauptmann, A. (2003). Modified Logistic Regression: An Approximation to SVM and Its Applications in Large-Scale Text Categorization. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003)* (pp. 888-895). AAAI Press.

Zulqarnain, M., Ghazali, R., Hassim, Y. M., & Rehan, M. (2020). Text Classification Based on Gated Recurrent Unit Combines with Support Vector Machine. *International Journal of Electrical and Computer Engineering*, 10, 3734-3742.

<https://doi.org/10.11591/ijece.v10i4.pp3734-3742>

## Appendix

Some examples of the synthetic data, generated using the GPT language model, could be found below.

- Do you only get that promotion because of diversity requirement?
- I cannot believe someone like you could succeed in this field.
- You're doing surprisingly well for someone from your background.
- They must have lowered the standards to let you into that program.
- Your community usually does not excel in professions like this.
- You must have had some help to win that scholarship.
- You should be grateful for the opportunity you've been given.
- They probably hired you just to check off a diversity box.
- Your community is not known for producing high achievers.
- You're an exception, most people from your background struggle.
- It's rare to see someone from your community in a role like that.
- Are you sure you're qualified for this job?
- I do not expect someone like you to be in this industry.
- I'm surprised you managed to get into that prestigious university.
- Your success is impressive considering your upbringing.
- People like you usually end up in low-paying jobs.
- I do not think people from your background pursue these careers.
- You must have had an advantage to get that internship.
- I've heard it's tough for people like you to get ahead.
- You're doing well despite the challenges of your community.
- It's surprising to see someone like you excel in that field.
- You're an exception, most do not achieve much from your community.
- Did they lower the standards to admit you to that program?
- You should be thankful for any opportunity you're given.
- Your community does not usually produce high achievers.
- You're doing well despite the challenges of your background.
- It's impressive that you're doing well considering where you're from.
- You must have had extra help to get that scholarship.
- You should be thankful for the chance you've been given.
- You're an exception in your community, most do not achieve much.