# CSSCUTN@DravidianLangTech: Abusive comments Detection in Tamil and Telugu

**Kathiravan Pannerselvam**[1], **Saranya Rajiakodi**[1], **Rahul Ponnusamy**[2], **Sajeetha Thavareesan**[3]

[1]Department of Computer Science, Central University of Tamil Nadu,
Thiruvarur, Tamil Nadu, India
[2]Insight SFI Research Centre for Data Analytics, University of Galway, Ireland,
[3]Eastern University, Sri Lanka.
saranya@cutn.ac.in

## Abstract

Code-mixing is a word or phrase-level act of interchanging two or more languages during a conversation or in written text within a sentence. This phenomenon is widespread on social media platforms, and understanding the underlying abusive comments in a code-mixed sentence is a complex challenge. We present our system in our submission for the DravidianLangTech Shared Task on Abusive Comment Detection in Tamil and Telugu. Our approach involves building a multiclass abusive detection model that recognizes 8 different labels. The provided samples are code-mixed Tamil-English text, where Tamil is represented in romanised form. We focused on the Multiclass classification subtask, and we leveraged Support Vector Machine (SVM), Random Forest (RF), and Logistic Regression (LR). Our method exhibited its effectiveness in the shared task by earning the ninth rank out of all competing systems for the classification of abusive comments in the code-mixed text. Our proposed classifier achieves an impressive accuracy of 0.99 and an F1-score of 0.99 for a balanced dataset using TF-IDF with SVM. It can be used effectively to detect abusive comments in Tamil, English code-mixed text.

## 1 Introduction

In the age of rapidly advancing technology and increased social media usage, online platforms have become integral to our daily lives, facilitating global communication and connection (Anita and Subalalitha, 2019; Thavareesan and Mahesan, 2019, 2020a,b; Subalalitha, 2019; Sakuntharaj and Mahesan, 2016, 2017, 2021). However, this interconnectedness has also led to the proliferation of abusive content and hate speech, posing severe challenges to maintaining a safe and respectful online environment (Bharathi and Agnusimmaculate Silvia, 2021; Bharathi and Varsha, 2022b; Swaminathan et al., 2022; Subramanian et al., 2022; Chinnaudayar Navaneethakrishnan et al., 2023). Automated classification of abusive comments is paramount to curb the spread of harmful content and protect users from online harassment (Priyadharshini et al., 2022). Moreover, abusive comments often emerge as a complex interplay of code-mixed text in multilingual communities, such as those in regions where Tamil and English speakers coexist and interact. Code-mixing is word or phrase level interchanging two or more languages within a single sentence or conversation (Chakravarthi et al., 2023a,b). The prevalence of code-mixed text in social media interactions further intensifies the difficulty of abusive comment classification, as traditional natural language processing techniques may not be directly applicable (Thara and Poornachandran, 2018; Chakravarthi et al., 2022a,b; ?). This research article addresses the challenges of abusive comment classification in Tamil and English code-mixed text using machine learning techniques and the Synthetic Minority Over-sampling Technique (SMOTE) to handle imbalanced data (Chawla et al., 2002; Kathiravan et al., 2023). We explore the effectiveness of employing SMOTE to create a balanced dataset, thereby mitigating the skewed distribution of abusive and non-abusive comments.

To achieve accurate and robust classification results, we explore two popular data representation techniques: Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF). BoW represents text data by counting the occurrence of words in documents, while TF-IDF considers the importance of words based on their frequency and rarity in the corpus (Akuma et al., 2022). We apply these representations to our balanced dataset and utilize three different machine learning algorithms: Support Vector Machines (SVM), Logistic Regression (LR), and Random Forest (RF).

One of the pivotal contributions of our research lies in showcasing the significance of SMOTE in enhancing the performance of the classifiers. By addressing the class imbalance issue, SMOTE enables the classifiers to better capture patterns and nuances related to abusive comments, improving classification accuracy and reducing the risk of misclassifying harmful content (Satriaji and Kusumaningrum, 2018). Our findings indicate that TF-IDF representation consistently outperforms BoW in SVM, LR, and RF models. The ability of TF-IDF to capture the semantic relevance and rareness of words provides a more substantial discriminatory power, leading to better classification results. Leveraging TF-IDF representation with SMOTE-balanced data yields an accuracy demonstrating the potential for more effective content moderation in code-mixed text environments.

However, despite the success achieved in our research, we recognize that addressing abusive comment classification in code-mixed text is an ongoing challenge. The dynamic nature of online content and the continuous evolution of abusive language demand continuous efforts in refining and updating our classification models.

## 2 Related works

The growing presence of Tamil and English, code-mixed data on social media platforms has sparked increased research interest in addressing abusive comments. As the multilingual nature of code-mixed content poses unique challenges, several studies have been conducted to develop effective methods for detecting and classifying abusive language in these contexts.

The research by Rajalakshmi et al. (2022) focuses on abusive comment detection in Tamil using multilingual transformer models. They explored the effectiveness of multilingual transformer models in handling the complexities of Tamil code-mixed text for abusive comment detection. The results demonstrated the potential of transformer-based approaches in achieving accurate and robust classification, significantly promoting a safer online environment. The research contributes valuable insights to the field of natural language processing for Dravidian languages, particularly in the context of abusive content moderation.

The research conducted by Prasanth et al. (2022) focuses on abusive comment detection in Tamil language text using the TF-IDF representation and the random kitchen sink algorithm. The authors addressed the challenge of identifying abusive content in the context of the Tamil language, which is particularly relevant for content moderation and user safety on social media platforms. The TF-IDF representation, known for capturing word importance and rarity, was combined with the random kitchen sink algorithm, a randomized feature mapping technique, to classify abusive comments.

Biradar and Saumya (2022) conducted research that focuses on the classification of abusive content in Dravidian code-mixed text using a Transformer-based approach. Their methodology leverages advanced natural language processing techniques to address the challenges of code-mixing in Tamil and other Dravidian languages. The study aims to improve content moderation systems and create a safer online environment by accurately detecting and handling abusive comments in multilingual communities. The paper provides valuable insights into using Transformers for abusive content classification in code-mixed text, contributing to the ongoing research efforts in Natural Language Processing (NLP) and content moderation.

Bharathi and Varsha (2022a) conducted the study on detecting abusive comments in the Tamil language using a Transformer-based approach. Transformers have proven to be highly effective in natural language processing tasks, and the authors explored their application for abusive comment classification in Tamil. Their work contributes to the growing body of research in content moderation and speech technologies for Dravidian languages, addressing the crucial issue of abusive language online. The detailed methodology, experimental setup, and results presented in the paper offer valuable insights into the effectiveness of the proposed approach, providing a significant contribution to the field of NLP for Tamil language processing.

The research conducted by Balouchzahi et al. (2022) focuses on abusive comment detection in the Tamil language. They employed a 1D Conv-LSTM model for classification. The study addressed the crucial issue of identifying abusive content in online Tamil text, considering the challenges posed by multilingual code-mixing. The proposed model demonstrated promising results in detecting abusive comments, and its effectiveness was showcased during the shared task. It also provides insights into the approach, methodology, and results achieved, contributing to developing content mod-

eration systems in multilingual environments.

## 3 System Description

In this section, we provide comprehensive information about the dataset and elaborate on the experiments conducted in our study. Figure 1 illustrates the system architecture designed for multiclass abusive comments classification using machine learning (ML) techniques, including SMOTE. The diagram depicts the overall flow of the classification process.
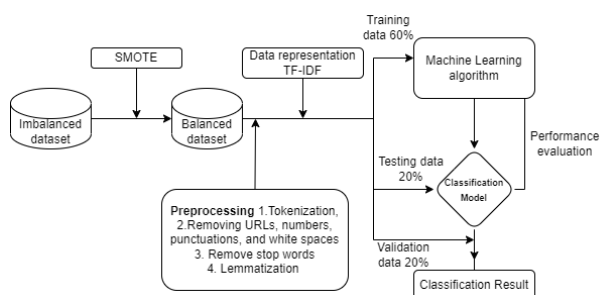


Figure 1: ML based Architecture.

### 3.1 Dataset

The dataset for the Abusive Comment Detection shared task, provided by the organizers, consists of code-mixed Tamil-English comments ("Priyadharshini et al., "2023") (Priyadharshini et al., 2022). It is important to note that the dataset is imbalanced, as depicted in Figure 2, showing the class-wise distribution of the data in percentages. The categories and their corresponding values are as follows: None-of-the-above (4633 comments), Misandry (1048 comments), Counter-speech (443 comments), Xenophobia (367 comments), Hope-Speech (266 comments), Misogyny (261 comments), Homophobia (215 comments), and Transphobic (197 comments). The imbalanced nature of the dataset poses challenges for effective classification and warrants careful consideration during model training and evaluation.

### 3.2 SMOTE

SMOTE (Synthetic Minority Over-sampling Technique) addresses the class imbalance in imbalanced datasets by generating synthetic samples for the minority class (Chawla et al., 2002). First, it identifies the minority class and selects k nearest neighbors for each instance within that class. Then, synthetic samples are created by interpolating between the
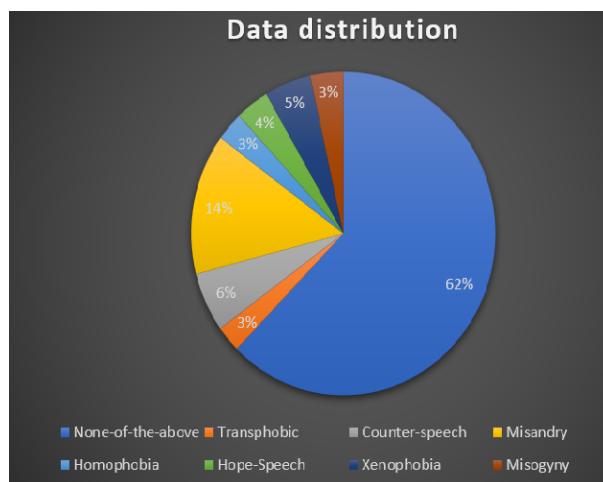


Figure 2: Data distribution

instance and one of its neighbors based on random proportions. This process is repeated for all instances in the minority class, resulting in a balanced dataset. The synthetic samples improve the representation of the minority class, allowing machine learning models to learn from a more diverse dataset and make more accurate predictions, particularly in abusive comment classification in code-mixed textual data (Shanmugavadivel et al., 2022). (Roy and Kumar, 2021) After applying the SMOTE to the initial imbalanced dataset, a significant transformation occurred, resulting in a balanced class distribution. Before SMOTE, the dataset had varying sample counts for each class, with the majority class "None-of-the-above" having 4633 comments. However, after SMOTE was applied, the number of samples in each class was equalized, with all classes containing 4633 comments. This process involved generating synthetic samples for the minority classes, effectively increasing their instances to match the majority class. The balanced dataset achieved through SMOTE addresses the class imbalance challenge, enabling machine learning models to train on a more representative and equitable dataset. Consequently, the classification models are better equipped to detect abusive comments in code-mixed Tamil-English data with enhanced accuracy and performance.

### 3.3 Preprocessing

We performed several preprocessing activities to clean and prepare the text data for further analysis. Firstly, we remove white spaces and punctuation to standardize the text and eliminate unnecessary things. Next, we eliminate stop words, which are

commonly occurring words with little semantic value, to reduce noise in the data. We then tokenize the text, breaking it into individual words or tokens, enabling more granular analysis (Kathiravan and Haridoss, 2018). Additionally, we apply lemmatization to reduce words to their base or root form, aiding in normalization and reducing word variations. These preprocessing activities collectively improve the quality of the text data, making it more suitable for subsequent steps in the abusive comment classification process using machine learning techniques.

## 3.4 Data representation

Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF) are standard text representation techniques used in natural language processing tasks. BoW converts a text document into a numerical vector by counting the frequency of each word present in the document (Qader et al., 2019). It disregards the word order and context, treating each word as an independent entity. The resulting vector represents the occurrence of words in the document, enabling the comparison and analysis of text data based on word frequencies. On the other hand, TF-IDF aims to capture the importance of words in a document relative to the entire corpus. It calculates the product of Term Frequency (TF), which is the frequency of a word in a document, and Inverse Document Frequency (IDF), which measures the rarity of the word across the entire corpus. TF-IDF assigns higher weights to words that are frequent in a document but rare in the corpus, indicating their significance in differentiating documents. As a result, TF-IDF representation allows for the emphasis of relevant terms while downplaying commonly occurring words, offering a more informative vector representation of the text data (Akuma et al., 2022). Both BoW and TF-IDF play crucial roles in text classification tasks, aiding in feature extraction and representing textual information in a format suitable for machine learning algorithms.

## 3.5 ML models

**Support Vector Machine (SVM)** is a robust and widely used supervised learning algorithm for classification and regression tasks. In classification, SVM aims to find the optimal hyperplane that best separates the data points of different classes in a high-dimensional space (Valero-Carreras et al., 2023) . The "support vectors"

are the data points closest to the hyperplane and play a crucial role in defining the decision boundary. SVM effectively handles both linearly separable and non-linearly separable data through kernel functions, such as polynomial or radial basis function (RBF) kernels, which map the data into a higher-dimensional space. SVM is known for handling high-dimensional data and generalizing to new, unseen data.

**Logistic Regression** is a popular statistical method used for binary classification tasks, where the goal is to predict the probability that a given data point belongs to a particular class (Xu et al., 2023). Despite its name, logistic Regression is primarily used for classification, not Regression. The algorithm models the relationship between the input features and the probability of the binary outcome using the logistic function (sigmoid function). The logistic function maps any real-valued number to the range [0, 1], which is then used to make the final classification decision. The model is trained by optimizing the parameters through maximum likelihood estimation. Logistic Regression is relatively simple, interpretable, and computationally efficient, making it a popular choice for binary classification tasks.

**Random Forest** is an ensemble learning method that combines multiple decision trees to achieve more accurate and robust predictions. Each decision tree is trained in a random forest on a random subset of the data (bootstrap samples) and a random subset of the features. This randomness helps reduce overfitting and increases the diversity among the individual trees (Das et al., 2023) . During prediction, the final output is determined by aggregating the predictions of all the trees, either through majority voting (for classification) or averaging (for Regression). Random forests are known for their ability to handle complex datasets, high-dimensional data, and non-linear relationships. They are also less prone to overfitting than a single decision tree and are widely used in machine learning tasks.

## 4 Result

In this research, we investigated the application of three machine learning algorithms, namely Support Vector Machine (SVM), Logistic Regression, and Random Forest, in conjunction with two data representation techniques: Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency

| Model | Feature set | Precision | recall | F1-score | Accuracy |
|---|---|---|---|---|---|
| Support Vector Machine | BoW | 0.98 | 0.97 | 0.97 | 0.97 |
| Random Forest | BoW | 0.98 | 0.98 | 0.98 | 0.98 |
| Logistic Regression | BoW | 0.98 | 0.97 | 0.97 | 0.97 |
| Support Vector Machine | TF-IDF | 0.99 | 0.99 | 0.99 | 0.99 |
| Random Forest | TF-IDF | 0.98 | 0.98 | 0.98 | 0.98 |
| Logistic Regression | TF-IDF | 0.97 | 0.97 | 0.97 | 0.97 |

Table 1: Evaluation metrics of various ML models.

(TF-IDF). After employing SMOTE to address the class imbalance, the models were evaluated for abusive comment classification in code-mixed Tamil-English text. This experimental research with TF-IDF representation and SMOTE-balanced data, SVM achieved an outstanding F1 score of 0.99, showcasing its ability to identify abusive comments in the dataset accurately. Logistic Regression, combined with TF-IDF, delivered an F1 score of 0.97, displaying competitive performance in abusive comment classification. Random Forest, with both BoW and TF-IDF representations on SMOTE-balanced data, attained F1 scores of 0.98 and 0.98, respectively. This highlights the algorithm's robustness in handling code-mixed text and class imbalance Table 1 illustrates it briefly.

## 5 Conclusion

In conclusion, this research delves into the critical issue of abusive comment classification in code-mixed text, particularly in Tamil and English code-mixed data. The escalating growth of social media and online platforms has amplified the need for compelling content moderation to ensure a safe and respectful digital environment. We have achieved promising results in accurately identifying abusive comments by harnessing the power of machine learning techniques and SMOTE for handling class imbalance. The utilization of SVM, Logistic Regression, and Random Forest in combination with TF-IDF representation has proven to be highly effective in capturing the nuances of abusive language, leading to improved classification accuracy. Notably, the application of SMOTE has significantly contributed to overcoming the challenges of class imbalance, enabling our models to make more reliable predictions.

**Limitation** is encountering out-of-vocabulary words or linguistic phenomena not accounted for during preprocessing. Code-mixing can introduce linguistic variations that may not be adequately handled by the existing language processing techniques, leading to potential misclassifications. Addressing these linguistic complexities in future research could enhance the model's performance and generalization capabilities. Notably, in the shared task, our approach secured the 9th rank out of all participating systems. Despite the competition's challenges and intense competition from other participants, our model's performance demonstrates its effectiveness in abusive comment classification in code-mixed text.

**Future work** investigating more advanced and context-aware feature extraction techniques will contribute to a more comprehensive analysis of code-mixed content. We continuously refine and update our models to develop robust and adaptive content moderation systems, creating a safer online environment for diverse, multilingual communities. we are optimistic that future work will address these limitations and elevate the performance of abusive comment classification in code-mixed text, contributing to a more inclusive and respectful online space for all users.

## References

Stephen Akuma, Tyosar Lubem, and Isaac Terngu Adom. 2022. Comparing Bag of Words and TF-IDF with different models for hate speech detection from live tweets. *International Journal of Information Technology*, 14(7):3629–3635.

R Anita and CN Subalalitha. 2019. An approach to cluster Tamil literatures using discourse connectives. In *2019 IEEE 1st International Conference on Energy, Systems and Information Processing (ICESIP)*, pages 1–4. IEEE.

Fazlourrahman Balouchzahi, Anusha Gowda, Hosahalli Shashirekha, and Grigori Sidorov. 2022. MUCIC@TamilNLP-ACL2022: Abusive Comment Detection in Tamil Language using 1D Conv-LSTM. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 64–69.

B Bharathi and A Agnusimmaculate Silvia. 2021. SS-NCSE_NLP@DravidianLangTech-EACL2021: Offensive language identification on multilingual code mixing text. In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 313–318, Kyiv. Association for Computational Linguistics.

B Bharathi and Josephine Varsha. 2022a. SSNCSE NLP@ TamilNLP-ACL2022: Transformer based approach for detection of abusive comment for Tamil language. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 158–164.

B Bharathi and Josephine Varsha. 2022b. SSNCSE NLP@TamilNLP-ACL2022: Transformer based approach for detection of abusive comment for Tamil language. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 158–164, Dublin, Ireland. Association for Computational Linguistics.

Shankar Biradar and Sunil Saumya. 2022. IIITDWD@TamilNLP-ACL2022: Transformer-based approach to classify abusive content in Dravidian code-mixed text. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 100–104, Dublin, Ireland. Association for Computational Linguistics.

Bharathi Raja Chakravarthi, Adeep Hande, Rahul Ponnusamy, Prasanna Kumar Kumaresan, and Ruba Priyadharshini. 2022a. How can we detect homophobia and transphobia? experiments in a multilingual code-mixed setting for social media governance. *International Journal of Information Management Data Insights*, 2(2):100119.

Bharathi Raja Chakravarthi, Manoj Balaji Jagadeeshan, Vasanth Palanikumar, and Ruba Priyadharshini. 2023a. Offensive language identification in Dravidian languages using MPNet and CNN. *International Journal of Information Management Data Insights*, 3(1):100151.

Bharathi Raja Chakravarthi, Ruba Priyadharshini, Shubanker Banerjee, Manoj Balaji Jagadeeshan, Prasanna Kumar Kumaresan, Rahul Ponnusamy, Sean Benhur, and John Philip McCrae. 2023b. Detecting abusive comments at a fine-grained level in a low-resource language. *Natural Language Processing Journal*, 3:100006.

Bharathi Raja Chakravarthi, Ruba Priyadharshini, Thenmozhi Durairaj, John McCrae, Paul Buitelaar, Prasanna Kumaresan, and Rahul Ponnusamy. 2022b. Overview of the shared task on homophobia and transphobia detection in social media comments. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 369–377, Dublin, Ireland. Association for Computational Linguistics.

Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.

Subalalitha Chinnaudayar Navaneethakrishnan, Bharathi Raja Chakravarthi, Kogilavani Shanmugavadivel, Malliga Subramanian, Prasanna Kumar Kumaresan, Bharathi, Lavanya Sambath Kumar, and Rahul Ponnusamy. 2023. Findings of shared task on sentiment analysis and homophobia detection of Youtube comments in code-mixed Dravidian languages. In *Proceedings of the 14th Annual Meeting of the Forum for Information Retrieval Evaluation*, FIRE '22, page 18–21, New York, NY, USA. Association for Computing Machinery.

Sunanda Das, Md Samir Imtiaz, Nieb Hasan Neom, Nazmul Siddique, and Hui Wang. 2023. A hybrid approach for bangla sign language recognition using deep transfer learning model with random forest classifier. *Expert Systems with Applications*, 213:118914.

P Kathiravan and N Haridoss. 2018. Preprocessing for mining the textual data-a review. *International Journal of Scientific Research in Computer Science Applications and Management Studies IJSRCSAMS*, 7(5).

P Kathiravan, P Shanmugavadivu, and R Saranya. 2023. Mitigating imbalanced data in online social networks using stratified k-means sampling. In *2023 8th International Conference on Business and Industrial Research (ICBIR)*, pages 883–888. IEEE.

SN Prasanth, R Aswin Raj, P Adhithan, B Premjith, and Soman Kp. 2022. CEN-Tamil@ DravidianLangTech-ACL2022: Abusive comment detection in Tamil using TF-IDF and random kitchen sink algorithm. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 70–74.

Ruba "Priyadharshini, Bharathi Raja Chakravarthi, Subalalitha Chinnaudayar Navaneethakrishnan, Malliga Subramanian, Kogilavani Shanmugavadivel, Premjith B, Abirami Murugappan, Sai Prashanth Karnati, Rishith, Chandu Janakiram, and Prasanna Kumar" Kumaresan. "2023". Findings of the shared task on Abusive Comment Detection in Tamil and Telugu. In *"Proceedings of the Third Workshop on Speech and Language Technologies for Dravidian Languages DravidianLangTech 2023"*. "Recent Advances in Natural Language Processing".

Ruba Priyadharshini, Bharathi Raja Chakravarthi, Subalalitha Cn, Thenmozhi Durairaj, Malliga Subramanian, Kogilavani Shanmugavadivel, Siddhanth U Hegde, and Prasanna Kumaresan. 2022. Overview of abusive comment detection in Tamil-ACL 2022. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 292–298, Dublin, Ireland. Association for Computational Linguistics.

311

Wisam A Qader, Musa M Ameen, and Bilal I Ahmed. 2019. An overview of bag of words; importance, implementation, applications, and challenges. In *2019 international engineering conference (IEC)*, pages 200–204. IEEE.

Ratnavel Rajalakshmi, Ankita Duraphe, and Antonette Shibani. 2022. DLRG@ DravidianLangTech-ACL2022: Abusive comment detection in Tamil using multilingual transformer models. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 207–213.

Pradeep Kumar Roy and Abhinav Kumar. 2021. Sentiment Analysis on Tamil Code-Mixed Text using Bi-LSTM. In *Working Notes of FIRE 2021-Forum for Information Retrieval Evaluation (Online). CEUR*.

Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2016. A novel hybrid approach to detect and correct spelling in Tamil text. In *2016 IEEE international conference on information and automation for sustainability (ICIAfS)*, pages 1–6. IEEE.

Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2017. Use of a novel hash-table for speeding-up suggestions for misspelt Tamil words. In *2017 IEEE international conference on industrial and information systems (ICIIS)*, pages 1–5. IEEE.

Ratnasingam Sakuntharaj and Sinnathamby Mahesan. 2021. Missing word detection and correction based on context of Tamil sentences using n-grams. In *2021 10th International Conference on Information and Automation for Sustainability (ICIAfS)*, pages 42–47. IEEE.

Widi Satriaji and Retno Kusumaningrum. 2018. Effect of synthetic minority oversampling technique (smote), feature representation, and classification algorithm on imbalanced sentiment analysis. In *2018 2nd International Conference on Informatics and Computational Sciences (ICICoS)*, pages 1–5. IEEE.

Kogilavani Shanmugavadivel, Sai Haritha Sampath, Pramod Nandhakumar, Prasath Mahalingam, Malliga Subramanian, Prasanna Kumar Kumaresan, and Ruba Priyadharshini. 2022. An analysis of machine learning models for sentiment analysis of Tamil code-mixed data. *Computer Speech & Language*, 76:101407.

CN Subalalitha. 2019. Information extraction framework for Kurunthogai. *Sādhanā*, 44(7):156.

Malliga Subramanian, Rahul Ponnusamy, Sean Benhur, Kogilavani Shanmugavadivel, Adhithiya Ganesan, Deepti Ravi, Gowtham Krishnan Shanmugasundaram, Ruba Priyadharshini, and Bharathi Raja Chakravarthi. 2022. Offensive language detection in Tamil youtube comments by adapters and cross-domain knowledge transfer. *Computer Speech & Language*, 76:101404.

Krithika Swaminathan, Bharathi B, Gayathri G L, and Hrishik Sampath. 2022. SSNCSE_NLP@LT-EDI-ACL2022: Homophobia/transphobia detection in multiple languages using SVM classifiers and BERT-based transformers. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 239–244, Dublin, Ireland. Association for Computational Linguistics.

S Thara and Prabaharan Poornachandran. 2018. Code-mixing: A brief survey. In *2018 International conference on advances in computing, communications and informatics (ICACCI)*, pages 2382–2388. IEEE.

Sajeetha Thavareesan and Sinnathamby Mahesan. 2019. Sentiment analysis in Tamil texts: A study on machine learning techniques and feature representation. In *2019 14th Conference on industrial and information systems (ICIIS)*, pages 320–325. IEEE.

Sajeetha Thavareesan and Sinnathamby Mahesan. 2020a. Sentiment lexicon expansion using Word2vec and fastText for sentiment prediction in Tamil texts. In *2020 Moratuwa engineering research conference (MERCon)*, pages 272–276. IEEE.

Sajeetha Thavareesan and Sinnathamby Mahesan. 2020b. Word embedding-based part of speech tagging in tamil texts. In *2020 IEEE 15th International conference on industrial and information systems (ICIIS)*, pages 478–482. IEEE.

Daniel Valero-Carreras, Javier Alcaraz, and Mercedes Landete. 2023. Comparing two svm models through different metrics based on the confusion matrix. *Computers & Operations Research*, 152:106131.

Yang Xu, Bern Klein, Genzhuang Li, and Bhushan Gopaluni. 2023. Evaluation of logistic regression and support vector machine approaches for xrf based particle sorting for a copper ore. *Minerals Engineering*, 192:108003.