

This document presents a study on unsupervised tweet categorization using a combination of statistical and semantic features, published in Multimedia Tools and Applications journal in 2023. The authors, Maibam Debina Devi and Navanath Saharia, aim to cluster tweets based on their content by utilizing term frequency-inverse document frequency (tf-idf) and a synonym-based weighting scheme. They employ a density-based clustering algorithm, specifically DBSCAN, with parameters set at minpoints=8 and epsilon=1.5, to categorize 1,000 tweets into six cohesive clusters, which are validated with a Silhouette coefficient score of 0.47.

The study introduces a domain-independent synonym-based semantic weighting scheme and develops a lightweight tweet categorization model that hybridizes statistical and semantic features. The authors also design a pre-processing sequence to preserve data integrity while removing noise.

The paper is structured into five sections, covering the introduction, literature survey, methodology, experiment and result analysis, and conclusion with future work directions.

The literature survey highlights the use of various techniques for tweet categorization, such as sentiment analysis, topic detection, and event categorization. The authors note a gap in domain-independent unsupervised tweet categorization and aim to address this with their approach, which prioritizes dataset-specific features.

The methodology section details the feature extraction process, including the use of tf-idf and the synonym-based feature extraction mechanism, and the clustering algorithm, with a focus on DBSCAN. The evaluation of the clustering algorithm is based on the Silhouette coefficient.

In the experiment and result analysis section, the authors describe the dataset of 1,000 English language hate-speech tweets and the pre-processing steps taken to normalize the data. They compare their approach with baseline methods like k-means and mini-batch k-means, demonstrating the superiority of their method in handling syntactic and semantic similarities.

The conclusion summarizes the key findings, emphasizing the effectiveness of the hybrid feature approach and the adaptability of the model for future enhancements. The authors acknowledge the financial support received and declare no conflicts of interest.

Overall, the study presents a novel approach to unsupervised tweet categorization that effectively utilizes semantic and statistical features, offering a domain-independent solution with potential for further research