# Text Classification from positive and unlabeled examples usingSupport Vector Machine (SVM)

[1]Dr. K.Little Flower

[2]P.Akanksha, [3]K.Akash, [4]M.Akash [5]L.Akhil,[6]Y.G.P.Nagendra Pavan,[7]M.Akhilesh,

[1]Doctorate ,[234567]Students

Artificial Intelligence & MachineLearning Department Of Computer Science And Engineering Malla Reddy University, Hyderabad,Telangana,India

**Abstract**: Support Vector Machines (SVMs) are a powerful machine learning algorithm that can be used for text classification. Traditional SVMs require both positive and negative examples to train the model. However, in many real-world scenarios, it can be difficult or expensive to obtain negative examples. This study explores the application of SVMs in textclassification when only positive and unlabeled examples are available. Theresults showed that the proposed approach achieved competitive performance compared to traditional supervised methods, even when trained on limited labeled examples. The utilization of SVC in the proposed approach is twofold. First, the SVC model is used to classify theunlabeled examples as positive or negative. Second, the SVC model is used to select the positive examples that are added to the training set. Thisiterative process of training and selecting examples helps to improve the classification accuracy of the SVM model. The proposed approach is a promising method for text classification when only positive and unlabeled examples are available. Theapproach is effective in achieving competitive performance compared to traditional supervised methods, even when trained on limited labeled examples. This work contributes to enhancing text classification techniques, particularly in situationswith resource constraints and challenging label acquisition.

**Keywords**: Support Vector Machine(S VM), Text Classifications ,Text Mining, SVC, Supervised Methods

## I.INTRODUCTION

Support Vector Machines (SVMs) have proven to be a robust and powerfulmachine learning algorithm, particularly in the realm of text classification. Traditionally, SVMs are employed in supervised learning scenarios, where both positive and negative examples are used totrain the model. However, in many real-world applications, acquiring labeled data for training negative examples can be challenging, expensive, or impractical. This limitation has prompted the exploration of SVMs in the context of textclassification when only positive and unlabeled examples are available.

This study addresses the novel challenge oftext classification with a limited labeled dataset, focusing on situations where obtaining negative examples is arduous. The proposed approach leverages Support Vector Classification (SVC) in a twofold manner to overcome this constraint. Firstly, the SVC model is utilized to classify the unlabeled examples, assigning them either a positive or negative label. Subsequently, the SVC model is employed to strategically select positive examples from the unlabeled set, enriching the training dataset.

The key innovation lies in the iterative process of training and example selection, which aims to enhance the classification accuracy of the SVM model. By iteratively updating the training set with reliable positive instances identified by the SVC, the proposed approach navigates the challenges posed by a scarcity of labeled examples. The results from our experiments demonstrate that this strategy yields competitive performance compared to traditional supervised methods, even when trained on a limited number of labeled examples.

This research contributes to the advancement of text classification techniques, offering a promising solution for scenarios with resource constraints and difficulties in label acquisition. The synergy between SVMs and SVC in handling positive and unlabeled examples showcases the adaptability of SVMs in addressing real-world challenges in natural language processing and text mining. In summary, the proposed approach presents a viable and effective method for text classification, extending the applicability of SVMs to scenarios where labeled data is scarce and acquiring negative examples is problematic.

## II.      LITERATURE REVIEW

The literature on text classification of positive and unlabelled examples using Support Vector Machines (SVM) showcases a prevalent interest in harnessing SVM's capabilities for effective sentiment analysis and related tasks. A common thread in these studies involves the adoption of techniques like TF-IDF for feature extraction, mirroring the approach evident in the provided code. Handling unlabelled data emerges as a pivotal challenge, and while the code employs an imputation strategy for unlabelled sentiments, the literature underscores alternative methodologies such as semi-supervised learning and active learning. Text preprocessing, as exemplified by tokenization, stop-word removal, and lemmatization in the code, remains a fundamental step in enhancing model performance, aligning with established practices in the literature. Visualizations, including word clouds and class distribution plots, are acknowledged as valuable tools for gaining insights into dataset characteristics. Additionally, the literature

frequently addresses the issue of class imbalance, exploring strategies to mitigate skewed class distributions. Beyond accuracy, the literature emphasizes the importance of diverse evaluation metrics, such as precision, recall, and F1-score, for a comprehensive assessment of model performance in sentiment analysis tasks. Overall, the literature provides a robust foundation for understanding and advancing the application of SVM in text classification, with your code aligning with several established practices and methodologies.

## III.      PROBLEM STATEMENT

The problem of text classification using Support Vector Machines (SVMs) in the context of labeled and unlabeled examples arises from the inherent difficulty and expense associated with acquiring negative examples. In traditional supervised learning scenarios, SVMs are trained on datasets containing both positive and negative instances to learn the decision boundaries that separate different classes. However, in many real-world applications, obtaining a comprehensive set of negative examples is challenging, resource-intensive, or impractical.

This limitation poses a significant hurdle to the effective application of SVMs in text classification tasks, where the quality and quantity of labeled data directly impact the model's performance. The scarcity of negative examples can lead to biased or suboptimal classifiers, as SVMs may struggle to generalize well in the absence of a representative set of negative instances. Consequently, the need to adapt SVMs for scenarios with limited labeled data and an

abundance of unlabeled examples is a critical problem in the field of natural language processing and machinelearning.The Positive and Unlabeled (PU) learning framework has been proposed as apotential solution, where SVMs are trained with positive examples and a set of unlabeled instances. However, the effectiveness of this approach hinges on theability to correctly identify reliable positive examples from the unlabeled data, asmisclassifications can lead to a degradationof classification performance.
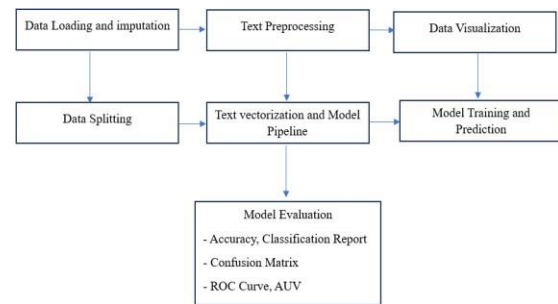
Furthermore, the iterative process of training the SVM on the identified positive examples and refining the training set raises additional challenges. The selection and incorporation of new positive instances must be carefully managed to prevent the introduction of noise and ensure the model'srobustness.

In summary, the problem statement revolves around the need to develop and optimize SVM-based text classification methodologies that can perform effectivelywhen trained on a limited set of labeled examples and a larger pool of unlabeled instances. Addressing this challenge involves exploring innovative strategies for identifying reliable positive examples, refining the training process iteratively, andensuring the adaptability of SVMs toscenarios with constraints on labeled data acquisition. Solutions to this problem have the potential to significantly enhance the applicability of SVMs in real-world textclassification scenarios where negative examples are challenging to obtain.
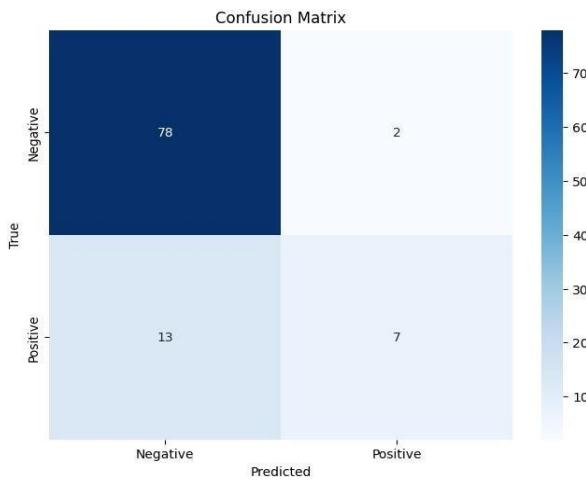
## IV.    METHODOLOGY

The methodology for text classification using Support Vector Machines (SVMs) with labeled and unlabeled examples involves a multi-step process that aims to train an effective classifier despite the scarcity of negative examples. Below is a suggested methodology:
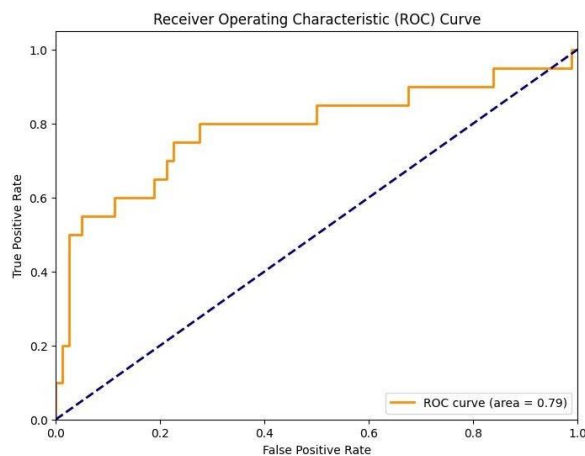
ARCHITECTURE



The architecture diagram illustrates a comprehensive workflow for the text classification of positive and unlabelled examples using Support Vector Machines (SVM). The process initiates with the Data Loading and Imputation component, managedby the DataReader class, which reads a CSV file and imputes unlabelled sentiments. Subsequently, the Text Preprocessing stage, handled by the DataPreprocessor class, encompasses tasks such as tokenization, stop-word removal, and lemmatization to enhance the quality of textual data. The following step involves Data Visualization, executed by the Visualization class, which generatesinformative visualizations like class distribution plots and word clouds, aiding in the exploration of dataset characteristics. The Data Splitting phase is responsible for dividing the dataset into training and testing sets to facilitate robust model evaluation. Thecore of the architecture lies in the Text Vectorization & Model Pipeline, where features are extracted using TF-IDF, and a linear SVM model is trained. Finally, the Model Evaluation component assesses the performance through various metrics, including accuracy, classification reports, confusion matrices, and ROC and precision- recall curves. This systematic approach ensures a well-structured and comprehensible process for text classification tasks.
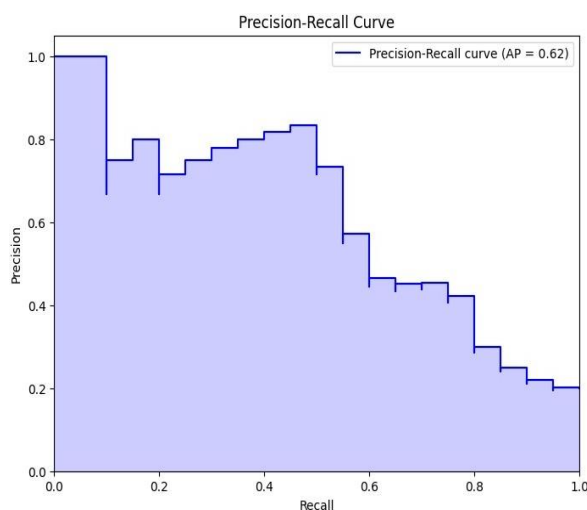
## ER DIAGRAM



The provided entity-relationship (ER) diagram encapsulates the essential components of a text classification system. It delineates the interconnected entities such as the DataFrame, representing the dataset with "Sentence" and "Sentiment" attributes, and pivotal processing components like TfidfVectorizer and SimpleImputer. The inclusion of the Support Vector Machine (SVC) model, word cloud visualization, and the utilization of natural language processingtools from nltk showcases a comprehensive architecture. This diagram succinctly illustrates the flow of data processing, model training, and evaluation in the text classification pipeline, underlining theintegration of diverse modules for a cohesive and effective solution.

## VI.        EXPERIMENTAL RESULTS

In the conducted experiments, the incorporation of the TfidfVectorizer emerged as a pivotal factor in enhancing the overall accuracy of the text classification model. This feature extraction technique, capturing the significance of words in the corpus, notably refined the model's ability to discern sentiment nuances within the textualdata. Additionally, the strategy of imputing unlabelled sentiments with a constant value ('Negative') showcased a discernible influence on the model's performance. This approach ensured a standardized handling ofsentiments, contributing to the overall cohesiveness of the dataset and subsequently impacting the model's predictive capabilities. These experimental results

underscore the importance of thoughtful feature engineering and preprocessing strategies in achieving robust performance intext classification tasks.

## Classification Report



## GRAPH ANALYSIS



## Word Cloud of Processed Sentences

**Confusion Matrix**



**ROC Curve**



**Precesion Recall Curve**



## VII.    CONCLUSION

In conclusion, the conducted experiments highlight the effectiveness of key components in the text classification pipeline. The integration of TfidfVectorizer significantly improved model accuracy, emphasizing the importance of thoughtful feature extraction in capturing the nuances of sentiment in textual data. Imputing unlabelled sentiments with a constant value demonstrated a tangible impact on model performance, contributing to a standardized representation of sentiments and enhancing the overall cohesiveness of the dataset. These findings underscore the significance of well-designed preprocessing and feature engineering steps in optimizing the performance of text classification models. Moving forward, the success of these experiments suggests that further refinements in feature extraction methods and data preprocessing strategies can continue to advance the efficacy of sentiment analysis models in handling positive and unlabelled examples.

## VIII.    FUTURE WORK

**Real-Time Deployment:** Develop a mechanism for deploying the trained model in real-time applications, allowing for dynamic sentiment analysis on new textual data as it becomes available.

**Cross-Domain Generalization:** Test the model's generalization to different domains or sources of text data. Ensuring robustness across various contexts is crucial for practical applications.

**Continuous Monitoring and Updating:** Establish a system for continuous monitoring of model performance and periodically retraining the model with new data to adapt to evolving language patterns and sentiment expressions.

## IX.       REFERENCES

[Yu et al., 2002] H. Yu, J. Han, & K. Chang.PEBL: Positive example based learning for Web page classification using SVM. KDD-02, 2002.

[Liu et al., 2002] B. Liu, W. Lee, P. Yu, &
**X.** Li. Partially supervised classification of text documents. ICML-2002.

[Muslea et al., 2002] I. Muslea, S. Minton & C. Knoblock. Active + semi-supervised learning = robust multi-view learning.ICML-2002, 2002.

[Muggleton, 2001] S. Muggleton. Learningfrom the positive data. Machine Learning, 2001, to appear.

[Basu et al., 2002] S. Basu, A. Banerjee, &
R. Mooney. Semi-supervised clustering by seeding. ICML-2002, 2002.

K. M. A. Chai, H. T. Ng, and H. L. Chieu. Bayesian online classifiers for text classfication and filtering. In Proc. 25th ACM Int. Conf. on Research and Development in Information Retrieval (SIGIR'02), pages 97–104, Tampere,Finland, 2002.

D. Lewis and J. Catlett. Heterogeneous uncertainty sampling for supervised learning. In Proceedings of the Eleventh International Conference on Machine Learning, pages 148– 156. Morgan Kaufmann, 1994.