**Paper Title:**
Article classification using natural language processing and machine learning

**Paper Link:**

**1 Summary**
**1.1 Motivation**

The paper is motivated by the need for efficient text classification in article submissions. With the rise of information and automatic computing, the authors aim to streamline the process by proposing an automated system for classifying submitted manuscripts into relevant topics, improving efficiency in the submission process.

**1.2 Contribution**

This paper significantly contributes to Natural Language Processing by introducing an automated system for efficient article classification. It streamlines the submission process through automated extraction of author information, title, and abstract, reducing manual effort. The study proposes effective data pre-processing techniques and leverages advanced methods like the vector space model and TF*IDF for accurate text representation. Additionally, the incorporation of a Vietnamese word segmentation tool enhances the system's accuracy, making it valuable for handling large volumes of article submissions with diverse topics.

**1.3 Methodology**

This paper introduces an automated system for article extraction and classification, tested on SVM, Naive Bayes, and kNN classifiers. The methodology involves a two-phase approach: training, where a model is generated using various machine learning algorithms and advanced techniques, and testing, where the model efficiently classifies articles in a testing dataset. The system streamlines the submission process, automating extraction and classification tasks with a focus on accuracy, especially in handling the Vietnamese language.

**1.4 Conclusion**
Among SVM, Naive Bayes, and KNN, Naive Bayes outperformed others. Specifically, Naive Bayes achieved 91.2%  average accuracy, surpassing previous models referenced in. In conclusion, Naive Bayes proved highly effective in classifying text articles.

**2 Limitations**
**2.1 First Limitation**

Algorithm Sensitivity: The effectiveness of the approach is contingent on the choice of machine learning algorithms, such as SVM, Naive Bayes, and kNN. Limitations inherent in these algorithms, including sensitivity to parameter tuning and potential biases in training data, may influence the system's overall performance and generalizability.

## 2.2 Second Limitation

Limited Generalization to Diverse Domains: The training dataset used in the study may be domain-specific, impacting the model's generalization capability to diverse topics. The system's effectiveness in classifying articles outside the scope of the training data, especially in broader or emerging domains, remains a potential limitation.

## 3 Synthesis

Despite language-dependence limitations and potential algorithm sensitivity, the methodology offers promising applications in broader linguistic contexts and diverse document types.

Future usage could involve extending the system to handle multiple languages, making it adaptable to diverse linguistic contexts. Additionally, the methodology could be explored in the context of other document types beyond articles, such as reports or essays, broadening its utility in various domains.