

## Chosen proposal text:

### Detection and Classification of Spam Using Machine Learning Algorithms

**Overall evaluation:** 1 (strong accept)

**Reviewer's confidence:** 3 (medium)

**Relevance of the contribution for an international computing education research audience:** 3 (appropriate and reasonably focused on computing education research, but limited in relevance to a specific locale, region, or country)

**Contribution to the computing education research field:** 4 (a clear contribution to the field)

**Discussion of related work:** 5 (all relevant work discussed and cited, and relationship to submission clearly described)

**Theoretical basis for the paper:** 3 (there is a theory there, but its relevance to the research is vague)

**Empirical basis for the paper:** 4 (data collected, good analysis)

**Writing and expression:** 4 (well written and expressed)

**Likelihood of generating discussion at the conference that will benefit the field:** 3 (neutral)

#### 1. Introduction

The reviewed proposal was mainly focused on differentiating spam or legitimate emails, by analysing the received emails with the machine learning technique. The proposal was overall acceptable with a clear point and a reasonable plan. The proposal scored 30 with the following reasons:

The proposal approach is evaluated using a publicly available dataset and compared against existing approaches. It is clearly presented about the goal, as well as what needs to be completed as a result. As pointed on the research topic, it highlights the clear and appropriate focus on computing education research, with the reason of necessitating spam detection methods. However, it did not specify clearly on the locale, region, or country; it was too vague and although spam content is a world-wise issue that cannot be limited into one area, it could be evaluated better. All related work were discussed and cited with the slide of sources, and the relationship to submission are all clearly described throughout the presentation. The literature reviews of the related works are presented with filtering techniques and classifying techniques; however, it was lacking with the explanation of the theoretical knowledge on each section. Especially during the methodology, since the techniques were not clearly demonstrated, it was too vague applying these techniques to the proposal aim. However, in terms of collecting and analysing the data, it has been well presented as it clearly shows the increase spam sites and demonstrates the approach to its issue. This includes presenting the smaller objectives, and main goal objectives separately to emphasise on their understanding of the chosen research topic. Especially during the methodology, the solution design with the solution design and implementation design provides the deep level of analysis, whilst demonstrating how the real-world dataset will be used to train the classifiers. However, it was vague on the step of acquiring a real-world dataset and the overall methodology seemed to be slightly overlapping with the project plan. As spam email is a current world-wise issue within the growth of the internet, this chosen topic is in a benefit of the current discussion at the conference as well as the solution of this problem will resolve few of the problems that are occurring in the current society. However, this is already on-going solving problem and already coming up with an idea of solutions, this might not be the best topic to be benefit in general.

## **2. Summary**

The presentation proposes a machine learning-based approach for detecting and classifying spam emails. It is mainly presenting about developing a machine learning algorithm to detect the emails' pattern and verify whether it is legitimate or spam email.

## **3. Methodology:**

The presented methodology is appropriate and of good quality. It shows the overall procedure of the plan, with the further explanations on the machine learning techniques and the results and conclusions. It also follows up with the solution design and implementation design, where it describes in detail about the algorithm section. However, the structure of the methodology seems like to be overlapping with the project plan, where it can be discussed fully about the steps of implementation rather than going through the entire procedure of the project. And due to this, it is lacking in detail on the specific hyperparameters used in the models and the rationale behind their selection.

## **4. Empirical Basis:**

The presentation demonstrates to be collecting and analysing data from a publicly available dataset to evaluate the proposed approach. And this dataset will then be used to train the machine learning algorithm, which is explained on the methodology part above. Especially during the solution design, it presents what types of machine learning will be implemented to build a classifier; Nearest Neighbour, Support Vector Machine and Naïve-Bayes, whilst explaining the concept of F-score to highlight how to perform the best categorization. Moreover, during the implementation design, it shows the brief code of how to put these algorithms into practice, such as merging multiple datasets, splitting into train, test and validation sets and the usage of sentiment analysis techniques. However, it seems to lack details on the dataset so could be better to provide more details on the dataset's characteristics and limitations and its solutions.

## **5. Use of Theory:**

The presentation provides enough theoretical basis for the proposed approach but could have linked better on its relevance to the research question. With the four literature reviews, it has briefly talked about the list of filtering algorithms, classifiers, spammer techniques and the methods of reducing the inaccuracy of the detection, however all literature reviews were not detailed enough to provide the enough pre-knowledge before getting into the research questions; it was mainly explaining about their functionality rather than the usage in the actual implementation.

## **6. Contributions/Results:**

As this presentation is just a framework of the actual plan, there was not many results that can be found, but so far it has been on the right track of the proposal that is being aimed for. It presents clear and well-analysed plan that demonstrate the effectiveness of the proposed approach in detecting and classifying the spam emails. It also provides a comparison with existing approaches, which further strengthens the contribution. Furthermore, the clear plan is well structured with the project plan sections, where it demonstrates on their procedure, by showing the audience, technology that will be used, mini goals and deliverables about the proposal. It again elaborates on the mini goal with the 4 steps of further exploration; Conduct, Design, Carry, and Analyse. Within these clear planning, it will potentially achieve the aim with high standard.

## **7. Significance:**

The finding of the proposal presentation is significant as they present an approach that can potentially improve the effectiveness of spam email detection and classification. The proposed approach can be used in various applications, including email filtering and even in cybersecurity areas. Therefore, this proposal can be useful to researchers, educators, and practitioners in the computing education and cybersecurity industry.