

# DETECTING SPAM E-MAIL WITH MACHINE LEARNING OPTIMIZED WITH BIO-INSPIRED METAHERUSTIC ALGORITMS.

Author<sup>1</sup>: SHAIK MASTAN,  
 Author<sup>2</sup>: THUMMA KARTHIK REDDY,  
 Author<sup>3</sup>: SHAIK NOORULA.

<sup>1</sup> Student, ECE, VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, A.P., INDIA

<sup>2</sup> Student, ECE, VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, A.P., INDIA

<sup>3</sup> Student, ECE, VASIREDDY VENKATADRI INSTITUTE OF TECHNOLOGY, A.P., INDIA

## ABSTRACT

*It investigates opinion mining by means of supervised learning techniques to search out the emotion of the student input bolstered characterized choices of instructing and learning. The examination led includes the apparatus of a blend of AI and common language preparing systems on understudy input data accumulated from module investigation overview consequences. Additionally, to offer a grade-by-grade clarification of the technique of accomplishment of opinion mining on or after scholar remarks using the open-source tool Python, the work additionally offers a comparative overall performance take a look remark supported, extracted alternatives like examination, teaching and so on. The consequences are as compared to be trying to find out higher overall performance with relevance several evaluation standards designed for the various techniques.*

**Keyword:** - Key word1: Support vector machine, Key word2: Genetic algorithm, Key word3: Coding language-python, Key word4: Designing -java script,html,css, key word5: particle swarm optimization.

## 1. INTRODUCTION

Machine Literacy models have been employed for multiple purposes in the field of computer wisdom from resolving a network business issue to detecting a malware. Emails are used regularly by numerous people for communication and for socializing. Security breaches that compromise client data allows spammers to imitate a compromised dispatch address to shoot illegitimate( spam) emails. This is also exploited to gain unauthorized access to their device by tricking the stoner into clicking the spam link within the spam dispatch, that constitutes a phishing attack. numerous tools and ways are offered by companies in order to descry spam emails in a network. Organizations have set up ltering mechanisms to descry unasked emails by setting up rules and conjuring the recall settings. Google is one of the top companies that offers 99.9 success in detecting similar emails. There are different areas for planting the spam lters similar as on the gate- way( router), on the pall hosted operations or on the stoner's computer. In order to overcome the discovery problem of spam emails, styles similar as content- grounded ltering, rule grounded ltering or Bayesian ltering have been applied. Unlike the knowledge engineering' where spam discovery rules are set up and are in constant need of homemade updating therefore consuming time and coffers, Machine literacy makes it easier because it learns to fete the unasked emails( spam) and licit emails( ham) automatically and also applies those learned instructions to unknown incoming emails. The proposed spam discovery to resolve the issue of the spam bracket problem can be farther experimented by point selection or automated parameter selection for the models. This exploration conducts trials involving v different machine literacy models with flyspeck mass Optimization( PSO) and inheritable Algorithm( GA). This will be compared with the base models to conclude whether the proposed models have bettered the performance with parameter tuning. The rest of this composition is organized as follows Section I presents the exploration to identify ways and styles used to resolve the bracket problem

## 2. RELATED WORK

### MACHINE LEARNING

Experimenters have taken a lead to apply machine literacy models to descry spam emails. In the paper, the authors have conducted trials with six different machine learning algorithms Naïve Bayes( NB) bracket, K- Nearest Neighbour( K- NN), Artificial Neural Network( ANN), Support Vector Machine( SVM), Artificial Immune System and Rough Sets. Their end of the trial was to imitate the detecting and recognising capability of humans. Tokenisation was explored and the conception handed two stages Training and Filtering. Their algorithm comported of four way DispatchPre-Processing, Description of the point, Spam Bracket and Performance Evaluation. It concluded that the Naïve Bayes handed the loftiest delicacy, perfection and recall.

### 3.PROPOSED WORK

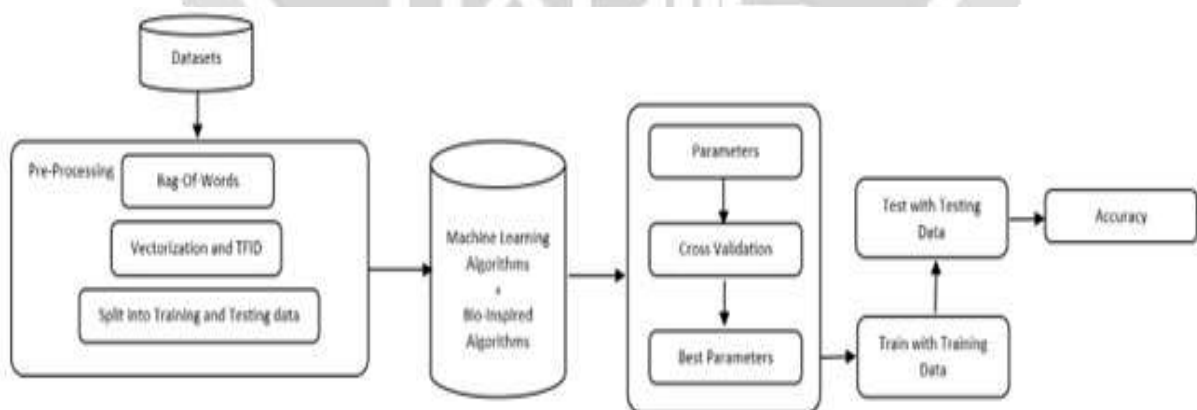
This exploration will experimentBio-inspired algorithms along with Machine literacy models. This will be conducted on different spam dispatch corpora that are intimately available. The paper aims to achieve the following objects

- 1) To explore machine literacy algorithms for the spam discovery problem.
- 2) To probe the workings of the algorithms with the acquired datasets.
- 3) To apply thebio-inspired algorithms.
- 4) To test and compare the delicacy of base models withbio-inspired perpetration.
- 5) To apply the frame using Python.

Scikit- Learn library will be explored to perform the trials with Python, and this will enable to edit the models, conductpre-processing and calculate the results. The program scripts will be enforced further with the optimization ways and compared with the base resultsi.e with dereliction parameters.

### 4.TOOLS AND TECHIQUES

WEKA is a GUI tool that allows to load a dataset and apply different functions/rules upon an algorithm. The application allows to apply the classification, regression, clustering algorithms and enable to visualise the data and the performance of the algorithm. An '.arff' file format of the spam datasets were fed into the program. Scikit-Learn (SKLearn) is an environment that is incorporated with Python programming language. The library offers a wide range of supervised algorithms that will be suitable for this project Keras is an API that supports Neural Networks. The API supports other deep learning algorithms for easy and fast approach. It offers CPU and GPU running capabilities in order to simultaneously process the models. Online tutorials are available for neural network for learning and development .



### 5.MODEL TRAINING AND TESTING PHASE

As discussed through the research, supervised learning methods were used and the model was trained with known data and tested with unknown data to predict the accuracy and other performance measures. To acquire the reliable results K-Fold cross validation was applied. This method does have its disadvantages such as, there is a chance that the testing data could be all spam emails, or the training set could include the majority

of spam emails. This was resolved by Stratified K-fold cross validation, which separates the data while making sure to have a good range of Spam and Ham into the distributed set.

## 6. MACHINE LAERING MODELS

The subsections below explain each of the Machine Learning models that will be implemented to achieve the aim of this work. The sections are accompanied with mathematical equations along with the pseudocode algorithms.

### 6.1 Naïve Bayes model

Naïve Bayes model is used to resolve bracket problems by using probability ways.

$$P(\text{WORD}|\text{Class}) \times P(\text{Class})$$

$$P(\text{Class}|\text{WORD}) = (P(\text{WORD}))$$

where WORD is( word1, word2,.. wordn) from within an uploaded dispatch and ' Class ' is either ' Spam ' or ' Ham '. The algorithm calculates the probability of a class from the bag of words handed by the program. Where  $P(\text{Class}|\text{WORD})$  is a posterior probability,  $P(\text{WORD}|\text{Class})$  is liability and  $P(\text{Class})$  is the previous probability.

### 6.2 Support Vector Machine

This algorithm plots each knot from a dataset within a dimensiona aeroplane and throughclassi\_cation fashion the cluster of data is separated by a hyperplane into their separate group. The hyperplane can be described as  $H = VX c$  where  $c$  is a constant and  $V$  is the vector. The SGDClassi\_er was loaded from scikit-learn library, which is the direct model with Stochastic Gradient Descent( SGD)', also known as the optimized interpretation of SVM. This algorithm provides more accurate results than SVM( SVC algorithm) itself. Disadvantage of working with SVC algorithm is that it can not handle a large dataset, whereas SGD provides efficiency and other tuning openings. The algorithm uses the literacy rate to reiterate over the sample data to optimize the Linear algorithm and it's denoted by the following equation  $1/\alpha(t)$  for the dereliction literacy rate as ' Optimal' where  $t$  is the time step which is acquired by multiplying number of duplications with number of samples( Emails).

### 6.3 Decision Tree Classifier

The Decision Tree model is grounded on the prophetic system. The model creates a order which is farther distributed intosub-categories and so on. The algorithm runs until the stoner has terminated or the program has reached its end decision. The model predicts the value of the data by learning from the handed training data. The longer and deeper the tree implies it has more complicated rules to beexecuted.Thepseudo-code for Decision Tree, where it terminates at the end of the knot for each split of the tree depth. analogous to MNB and SGD, Decision Tree( DT) algorithm was loaded from the Scikit- learn library and it's executed on the dereliction parameters which are Gini' for Criterion and ' stylish ' for Splitter. The advantage of Gini is that it calculates the inaply labelled data that was named aimlessly.

### 6.4 Random Forest Classifier

Random Forest (RF) algorithm can be used for both classification and regression. The algorithm predicts the classes by using decision tree, where each tree predicts the classification class. This is evaluated by the RF model to select the high number of predicted class as an assigned prediction.

## 7.INPUTS AND OUTPUTS

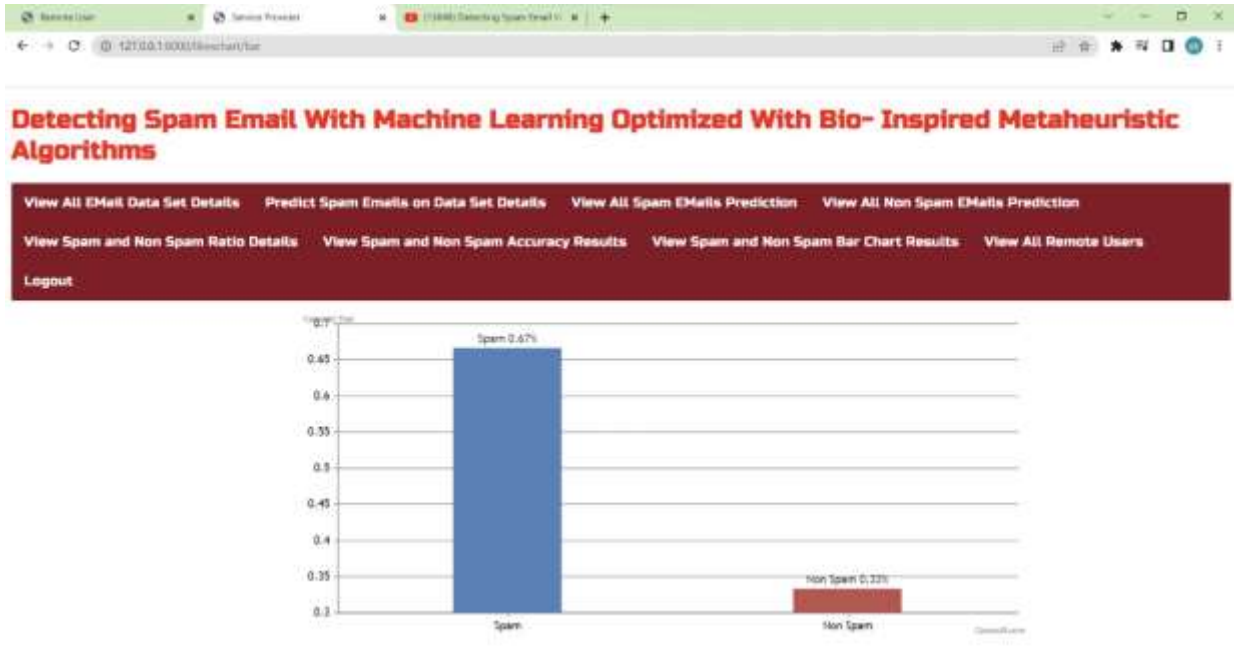


Fig : Bar chart of the spam and non-sapm e-mails

The screenshot shows a search interface with a search bar and a table of email details. The search bar contains the text 'Lotteries are always attracts'. The table below has the following data:

From	Subject	To	Email Date	Mailed_by	Signed_by	Security	Contents
shobh123@gmail.com	favour of lotteries	Aaniya123@gmail.com	18/10/2020	aa33.benkartz.in	shobh123@gmail.com	Standard encryption (TLS)	Lucky draw and lotteries are always attractive and amazing. So it always makes us a little happier, sometimes giving all the information or paying for the Lucky Draw Complaint Contest and lottery. But it

**8. BIO-INSPIRED OPTIMIZATION ALGORITHMS**

The PSO is grounded on the swarming styles observed in, sh or catcalls. The patches are estimated grounded on their stylish position and overall global position. patches within a hunt space are scattered toa and the global stylish position. The Pyswarms library offers different computations and ways for PSO to be used with an ML model similar as point subset selection or parameter tuning optimization. As delved in the former sections, the point selection can reduce point space but can also discard some features that can be useful during the bracket. thus, PSO will be used to tune and the hyperactive- parameter for a given ML/ NN model.

**9. CONCLUSION**

The project successfully implemented models combined with bio-inspired algorithms. The spam email corpus used within the project were both numerical as well as alphabetical. Approximately 50,000 emails were



tested with the proposed models. The numerical corpuses (PU), had restrictions in terms of feature extraction as the words were replaced by numbers. But the alphabetical corpuses performed better in terms of extraction of the features and predicting the outcome. Initially, WEKA acted as a black box that ran the datasets on 14 different classification algorithms and provided the top 4 algorithms: Multinomial Naïve Bayes, Support Vector Machine, Random Forest and Decision Tree. These algorithms were then tested and experimented with Scikit-learns library and its modules. This resulted in upgrading the SVM module with SGD classifier, which acts the same as SVM but performs better on the large datasets. SGD was implemented using Python and experimented with feature extraction and stop words removal along with converting the tokens for the algorithms to process. Genetic Algorithm worked better overall for both text-based datasets and numerical-based datasets than PSO. The PSO worked well for Multinomial Naïve Bayes and Stochastic Gradient Descent, whereas GA worked well for Random Forest and Decision Tree. Naïve Bayes algorithm was proved to have been the best algorithm for spam detection. This was concluded by evaluating the results for both numerical and alphabetical based dataset. The highest accuracy provided was 100% with GA optimization on randomized data distribution for 80:20 train and test split set on Spam Assassin dataset. In terms of F1-Score, precision and recall, Genetic Algorithm had more impact PSO on MNB ,SGD,RF and DT.

## 10. FUTURE WORK

We plan to further carry out the machine learning algorithms to optimize and compare with different bio-inspired algorithms such as Fire\_y, Bee Colony and Ant Colony Optimization as researched in the previous sections. We could also explore the Deep learning Neural Network with PSO and GA by exploring different libraries such as TensorFlow's DNN Classifier or similar. We found that the Neural Network algorithm could have worked better with more dimension like providing broader range of values for learning rate, activation, solver, and alpha. If this project is taken further, implementation for MLP could be done through Keras or TensorFlow with GPU application. This will allow the user to input other parameters and a range of possibilities as their key values. The user can consider implementing the PSO objective Function with RSCV to compare the difference for accuracy improvement. The PSO and GA can provide better accuracies by incorporating NLP techniques.

## 11. REFERENCES

1. W. Feng, J. Sun, L. Zhang, C. Cao, and Q. Yang, "A support vector machine based Naive Bayes algorithm for spam filtering," in Proc. IEEE 35th Int. Perform. Comput. Commun. Conf. (IPCCC), Dec. 2016, pp. 18, Doi: 10.1109/pccc.2016.7820655.
2. E. G. Dada, J. S. Bassi, H. Chiroma, S. M. Abdulhamid, A. O. Adetunmbi, and O. E. Ajibuwa, "Machine learning for email spam filtering: Review, approaches and open research problems," Heliyon, vol. 5, no. 6, Jun. 2019, Art. no. e01802, Doi: 10.1016/j.heliyon. 2019.e01802.
3. K. Agarwal and T. Kumar, "Email spam detection using integrated approach of Naïve Bayes and particle swarm optimization," in Proc. 2nd Int. Conf. Intell. Comput. Control Syst. (ICICCS), Jun. 2018, pp. 685690, doi: 10.1109/ICCONS.2018.8662957.
4. A. I. Taloba and S. S. I. Ismail, "An intelligent hybrid technique of decision tree and genetic algorithm for E-Mail spam detection," in Proc. 9th Int. Conf. Intell. Comput. Inf. Syst. (ICICIS), Cairo, Egypt, Dec. 2019, pp. 99104, Doi: 10.1109/ICICIS46948.2019.9014756.
5. S. Mohammed, O. Mohammed, and J. Fiaidhi, "Classifying unsolicited bulk email (UBE) using Python machine learning techniques," Int.J. Hybrid Inf. Technol., vol. 6, no. 1, pp. 4355, 2013. [Online]. Available: [https://www.researchgate.net/publication/236970412\\_Classifying\\_Unsolicited\\_Bulk\\_Email\\_UBE\\_using\\_Python\\_Machine\\_Learning\\_Techniques](https://www.researchgate.net/publication/236970412_Classifying_Unsolicited_Bulk_Email_UBE_using_Python_Machine_Learning_Techniques).
6. A. Wijaya and A. Bisri, "Hybrid decision tree and logistic regression classifier for email spam detection," in Proc. 8th Int. Conf. Inf. Technol. Electr. Eng. (ICITEE), Oct. 2016, pp. 14, Doi: 10.1109/ICITEED.2016.7863267.