

Journal of Advances in Mathematical & Computational Sciences
An International Pan-African Multidisciplinary Journal of the SMART Research Group
International Centre for IT & Development (ICITD) USA
© Creative Research Publishers
Available online at <https://www.isteams.net/mathematics-computationaljournal.info>
CrossREF Member Listing - <https://www.crossref.org/06members/50go-live.html>

Hybrid Genetic Algorithm Trained Bayesian Ensemble for Short Messages Spam Detection

¹Odiakaose Christopher, ² Emordi Frances, ³Ejeh Patrick, ⁴Ashioba Nwanze, ⁵Odeh Christopher, ⁵Attoh Obiageli & ⁶Azaka Maduabuchuku

^{1,3,4,5,6} Dept of Computer Science, Dennis Osadebay University Anwai-Asaba, Nigeria.

^{2,5}Department of Cybersecurity, Dennis Osadebay University Anwai-Asaba, Nigeria.

E-mails: osegalaxy@gmail.com, emordi.frances@dou.edu.ng, patrick.ejeh@dou.edu.ng, ashioba.nwanze@dou.edu.ng, odeh.christopher@dou.edu.ng, Attoh,obiageli@dou.edu.ng, azaka.maduabuchuku@dou.edu.ng

ABSTRACT

Today's popularity of the short messages services (SMS) has created a propitious environment for spamming to thrive. Spams are unsolicited advertising, adult-themed or inappropriate content, premium fraud, smishing and malware. They are a constant reminder of the need for an effective spam filter. However, SMS limitations of 160-characters and 140-bytes size as well as its being riddled with slangs, emoticons and abbreviations further inhibits effective training of models to aid accurate classification. The study proposes Genetic Algorithm Trained Bayesian Network solution that seeks to normalize noisy feats, expand text via use of lexicographic and semantic dictionaries that uses word sense disambiguation technique to train the underlying learning heuristics. And in turn, effectively help to classify SMS in spam and legitimate classes. Hybrid model comprises of text preprocessing, feature selection as well as training and classification section. Study uses a hybrid Genetic Algorithm trained Bayesian model for which the GA is used for feature selection; while, the Bayesian algorithm is used as classifier.

Keywords: Hybrid Genetic Algorithm, Trained Bayesian Ensemble, Short Messages Spam Detection

Odiakaose, C. Emordi, F. Ejeh, P., Ashioba, N., Odeh, C., Attoh, O. & Azaka, M. (2024): Hybrid Genetic Algorithm Trained Bayesian Ensemble for Short Messages Spam Detection. *Journal of Advances in Mathematical & Computational Science*. Vol. 12, No. 1. Pp 37-52.
Available online at www.isteams.net/mathematics-computationaljournal. dx.doi.org/10.22624/AIMS/MATHS/V12N1P4

1. INTRODUCTION

The advent of smartphones with enhance features has contributed to huge adoption of short messaging by users due to its portability, mobility, ubiquity of services and its low cost continues to promote SMSM to become the most used means of communication globally. ,



The most common challenge consists of distorted images and text. To pass the challenge, a human must type the text or arrange the images correctly. With challenge/response false positives can be reduced to barest minimum. Another merit of this approach is in its low system resource requirements, since no CPU-intensive pattern matching is required. However, this approach causes more problems than it solves. For inexperienced or visual handicapped users, the challenges are completely unsolvable. Regular users are provoked by the challenges and choose not to do so since they view it as an unacceptable irritation. Also, automated email that a user would want to receive (travel confirmations, online purchase receipts, etc) are trapped by this approach and never delivered (Process Software report, 2006).

2.2. Motivation / Statement of Problem

Study is motivated (Ifeka & Akinbobola, 2015; Igwenagu, 2015) as thus:

1. Spams have continued to soar with the advent of SMS. The alarming growth rate of spams with SMS popularity have now created a propitious environ for spammers to exploit subscribers; Thus, causing both financial loss and emotional instability as consequences to users, corporate organs and mobile network operator(s).
2. Academic researches and companies are today, faced with the challenge of dealing with SMS spam. A major issue has been that existing approach(es) to resolving SMS spam are imported from successful email anti-spam solutions (Wang *et al.*, 2010). Thus, are quite unable to effectively and efficiently tackle SMS spam successfully – as their performance is seriously hampered and degraded by the parametric feats used to filter spams.
3. The formulation and design of an effective SMS filter has continued to suffered setback(s) due to the inherent reason that SMS filters by design are not as simple as email filters due to its limited size of 160-characters of 140bytes sized data. These amongst other constraints, continue to create rippled impediment in size of feature to be selected for training and consequently contributing to poor learning and classification of learning algorithm.
4. Also, SMS is rippled with abbreviations, slangs and emoticons that inhibit proper classification of words and/or text corpus (Tiago *et al.*, 2016).

To overcome these, we deploy a hybrid SMS filter ensemble to reduce noise in form of slangs, emoticons, abbreviations as well as expand message size to enhance effective classification using text normalization and semantic expansion in SMS spam filtering.

2.3. Data Sampling

Dataset is retrieved from the Knowledge in discovery dataset (KDD-CUP 1999). Also, the dataset is split into: **training** (75%), and **test** (25%). (Ojugo & Eboka, 2021; Ojugo & Ekurume, 2021b, 2021a; Ojugo & Nwankwo, 2021b).

2.4. The Proposed Memetic Ensemble

GA is inspired by evolution, and consist of population based on potential solutions to a specific task. An individual with gene close to optimal, is fit. Fitness function determines how close an individual is to optimal solution. Each solution is an individual whose optimal is found via 4-operators (Mohd Ibrahim *et al.*, 2022; Tomar & Manjhar, 2015; Voke *et al.*, 2023).

2.5. The Experimental Framework

Figure 1 shows the schematics of the proposed experimental model (Ojugo & Eboka, 2021; Ojugo & Ekurume, 2021b, 2021a; Ojugo & Nwankwo, 2021b), which is explained in sections as:

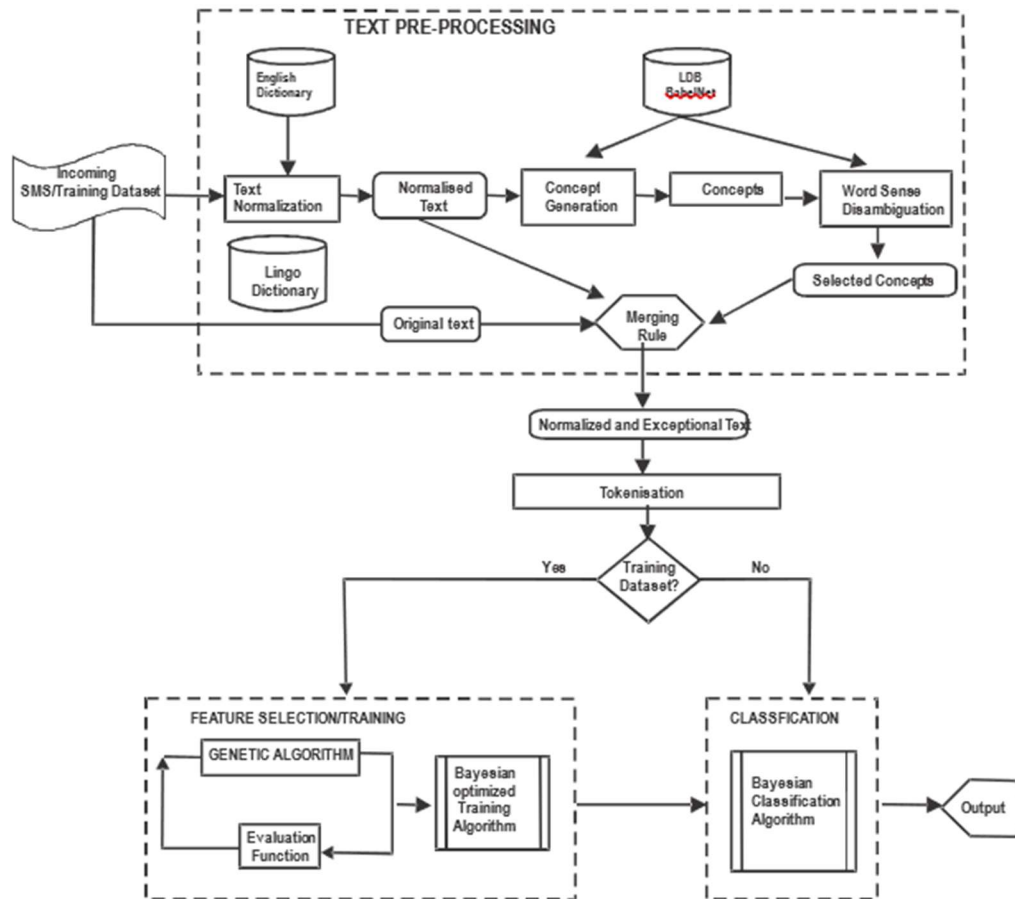
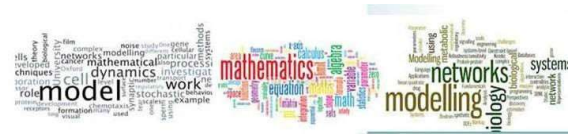


Figure 1. Proposed Memetic Model (Ojugo, Yoro, et al., 2013; Ojugo & Yoro, 2013, 2020, 2021)

Pre-processing Section: is accomplished as:

1. Raw text refers to the original text from the sender for normalization and expansion.
2. Text normalization uses two dictionaries. The first is an English dictionary to check if the text are English so as to then normalize text to its root form. The second is a slang dictionary used to translate slangs into English text. The basic operation of this stage is to replace slangs and abbreviation with standard English words from this dictionary. The Freeling English dictionary and No slang dictionary are proposed.
3. Concepts generation are semantically analyzed already normalized text to deduce their concept. The concepts are provided by Language Data Base BabelNet repository.



- Ojugo, A. A., & Eboka, A. O. (2018a). Assessing Users Satisfaction and Experience on Academic Websites: A Case of Selected Nigerian Universities Websites. *International Journal of Information Technology and Computer Science*, 10(10), 53–61. <https://doi.org/10.5815/ijitcs.2018.10.07>
- Ojugo, A. A., & Eboka, A. O. (2018b). Comparative Evaluation for High Intelligent Performance Adaptive Model for Spam Phishing Detection. *Digital Technologies*, 3(1), 9–15. <https://doi.org/10.12691/dt-3-1-2>
- Ojugo, A. A., & Eboka, A. O. (2018c). Modeling the Computational Solution of Market Basket Associative Rule Mining Approaches Using Deep Neural Network. *Digital Technologies*, 3(1), 1–8. <https://doi.org/10.12691/dt-3-1-1>
- Ojugo, A. A., & Eboka, A. O. (2020a). An Empirical Evaluation On Comparative Machine Learning Techniques For Detection of The Distributed Denial of Service (DDoS) Attacks. *Journal of Applied Science, Engineering, Technology, and Education*, 2(1), 18–27. <https://doi.org/10.35877/454ri.asci2192>
- Ojugo, A. A., & Eboka, A. O. (2020b). Memetic algorithm for short messaging service spam filter using text normalization and semantic approach. *International Journal of Informatics and Communication Technology (IJ-ICT)*, 9(1), 9. <https://doi.org/10.11591/ijict.v9i1.pp9-18>
- Ojugo, A. A., & Eboka, A. O. (2021). Empirical Bayesian network to improve service delivery and performance dependability on a campus network. *IAES International Journal of Artificial Intelligence (IJ-AI)*, 10(3), 623. <https://doi.org/10.11591/ijai.v10.i3.pp623-635>
- Ojugo, A. A., Eboka, A. O., Yoro, R. E., Yerokun, M. O., & Efozia, F. N. (2015). Hybrid Model for Early Diabetes Diagnosis. *2015 Second International Conference on Mathematics and Computers in Sciences and in Industry (MCSI)*, 50, 55–65. <https://doi.org/10.1109/MCSI.2015.35>
- Ojugo, A. A., & Ekurume, E. O. (2021a). Deep Learning Network Anomaly-Based Intrusion Detection Ensemble For Predictive Intelligence To Curb Malicious Connections: An Empirical Evidence. *International Journal of Advanced Trends in Computer Science and Engineering*, 10(3), 2090–2102. <https://doi.org/10.30534/ijatcse/2021/851032021>
- Ojugo, A. A., & Ekurume, E. O. (2021b). Predictive Intelligent Decision Support Model in Forecasting of the Diabetes Pandemic Using a Reinforcement Deep Learning Approach. *International Journal of Education and Management Engineering*, 11(2), 40–48. <https://doi.org/10.5815/ijeme.2021.02.05>
- Ojugo, A. A., & Nwankwo, O. (2021a). Forging a Spectral-Clustering Multi-Agent Hybrid Deep Learning Model To Predict Rainfall Runoff In Nigeria. *International Journal of Innovative Science, Engineering and Technology*, 8(3), 140–147.
- Ojugo, A. A., & Nwankwo, O. (2021b). Spectral-Cluster Solution For Credit-Card Fraud Detection Using A Genetic Algorithm Trained Modular Deep Learning Neural Network. *JINAV: Journal of Information and Visualization*, 2(1), 15–24. <https://doi.org/10.35877/454RI.jinav274>
- Ojugo, A. A., Odiakaose, C. C., Emordi, F. U., Ejeh, P. O., Adigwe, W., Anazia, K. E., & Nwozor, B. (2023). Forging a learner-centric blended-learning framework via an adaptive content-based architecture. *Science in Information Technology Letters*, 4(1), 40–53. <https://doi.org/10.31763/sitech.v4i1.1186>
- Ojugo, A. A., & Okobah, I. P. (2017a). Computational Solution for Modeling Rainfall Runoff Using Intelligent Stochastic Model: A Case of Warri in Delta State Nigeria. *Journal of Digital Innovations and Contemporary Res. in Science Engineering and Technology*, 5(4), 45–58. <https://doi.org/10.22624>

