



Quarterly Journal of **Optimization In Soft Computing**

Vol. 2, Issue 3, Autumn 2024

- **Bridging Technology and Language: Exploring Soft Computing Solutions for Effective English Language Teaching in Iran**
Ali Jahanbakhsh and Mehdi Jahangiri
- **A Novel Approach for Intrusion Detection System in IoT Using Correlation-Based Hybrid Feature Selection and Harris Hawk Optimization Algorithm**
Yashar Salami, Yaser Ebazadeh, Mehdi Hamrang and Nooshin Allahbakhshi
- **Improving the Performance of Permanent Magnet Synchronous Motor by Using Direct Current Control Method Based On Predictive Controller with Continuous Control Set**
Hamid Rahimi Esfahani and Reza Sharifian Dastjerdi
- **A Feature Selection Method on Gene Expression Microarray Data for Cancer Classification Abstract**
Parham Kiyoumars, Farshad Kiyoumars, Behzad Zamani and Mohammad Karbasiyoun
- **Comparison of Optimal Artificial Neural Network Models for Groundwater Nitrate Simulation (Case Study: Behbahan Plain)**
Atefeh Sayadi Shahraki, Fahimeh Sayadi Shahraki and Bijan Haghigati
- **Transformer-based Meme-sensitive Cross-modal Sentiment Analysis Using Visual-Textual Data in Social Media**
Zahra Pakdaman, Abbas Koochari and Arash Sharifi



Paper Type (Research paper)

Bridging Technology and Language: Exploring Soft Computing Solutions for Effective English Language Teaching in Iran

Ali Jahanbakhsh¹, Mehdi Jahangiri^{2,*}

1. College of Skills and Entrepreneurship, Shahrekord Branch, Islamic Azad University, Shahrekord, Iran

2. Energy Research Center, Shahrekord Branch, Islamic Azad University, Shahrekord, Iran.

Article Info

Article History:

Received: 2024/08/08

Revised: 2024/08/20

Accepted: 2024/10/05

DOI:

Keywords:

Soft Computing; English Language Teaching; Optimization Algorithms; Learning-Oriented Assessment; Fuzzy logic; Artificial Neural Networks.

*Corresponding Author's Email Address:

Mehdi_Jahangiri@yahoo.com

Abstract

The rapid evolution of technology and the increasing complexity of educational environments necessitate innovative approaches to language instruction. This paper explores the intersection of optimization techniques in soft computing and their application to English language teaching (ELT) in Iran. The necessity and importance of this review stem from the challenges faced by educators in adapting traditional teaching methodologies to meet the diverse needs of learners in a rapidly changing digital landscape. This study highlights their potential to enhance personalized learning experiences, improve curriculum design, and facilitate adaptive assessment strategies by synthesising existing literature on soft computing methods such as fuzzy logic, neural networks, and genetic algorithms. The work innovatively integrates these optimization techniques into ELT frameworks, proposing a model that leverages data-driven insights to tailor instructional strategies according to individual learner profiles. Key findings reveal significant improvements in student engagement, retention rates, and language proficiency when soft computing methods are employed. Moreover, the results indicate that such approaches can address the unique linguistic and cultural challenges faced by Iranian learners, fostering a more inclusive and effective educational environment. This paper contributes to the ongoing discourse on technology-enhanced language education by providing evidence of the benefits of optimization in soft computing. It underscores the imperative for educators and policymakers in Iran to embrace these methodologies to transform ELT, ultimately equipping learners with the skills necessary to thrive in an interconnected world.

1. Introduction

Optimization in Soft Computing and its relation with ELT in Iran examines the interdisciplinary application of advanced computational techniques to enhance educational outcomes in English language instruction [1]. Soft computing, an umbrella term encompassing methodologies like Fuzzy Logic (FL), Evolutionary Computation (EC), Machine Learning (ML), and Probabilistic Reasoning (PR), offers flexible and approximate solutions to complex, high-level problems in contrast to rigid traditional computing approaches

[2]. This versatility makes soft computing particularly valuable in fields requiring nuanced decision-making and problem-solving, including education [3].

The integration of soft computing in ELT in Iran is notable for its potential to transform traditional pedagogical practices [4]. English proficiency plays a crucial role in the effective use of technology in education, as it helps teachers and students navigate and utilize digital tools [5]. By applying optimization algorithms, educators can

create more personalized and efficient learning environments that cater to the specific needs of students. For instance, the use of Learning-Oriented Assessment (LOA) has shown to improve language acquisition by aligning continuous assessment with learning goals, thereby fostering a more engaging and contextualized learning experience [6].

Despite its benefits, the adoption of soft computing techniques in ELT faces challenges, such as time constraints, large class sizes, and a lack of proper training for teachers [6]. Additionally, the prevailing exam-oriented culture can conflict with the formative and continuous nature of LOA practices, posing further obstacles to its widespread implementation [6]. Addressing these issues requires a concerted effort to provide adequate training and resources for educators, along with a shift towards more holistic and integrative assessment methods.

The future prospects of applying optimization in soft computing to ELT in Iran are promising, with

potential advancements in data mining algorithms and deep learning techniques like Non-Convex Particle Swarm Optimization (NC-PSO) combined with Generative Adversarial Networks (GANs) offering new avenues for personalized instruction [7]. By leveraging these technologies, educators can optimize teaching strategies and enhance the overall quality of English education, ultimately fostering a more adaptive and effective learning environment.

The novelty of this work lies in its potential to enhance the quality of English language education in Iran and bridge the gap between technology and language. Additionally, by utilizing the proposed solutions, the learning experience for language learners can be made richer and more effective.

The flowchart of the current work is displayed in Figure 1. Outlining the various sections and subsections aids in a better understanding of the article and helps to solidify its structure in the reader's mind.

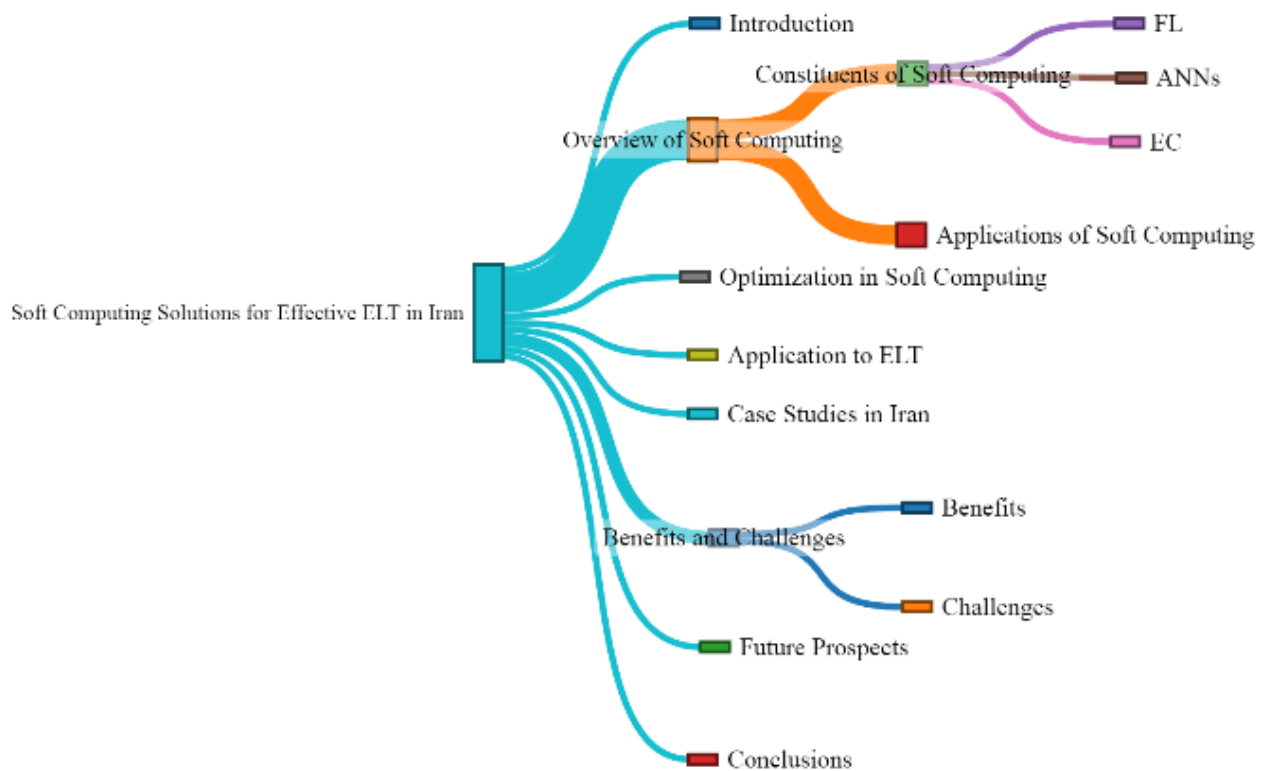


Figure 1. Flowchart of present study

2. Overview of Soft Computing

Soft computing is an umbrella term used to describe types of algorithms that produce approximate solutions to high-level, unsolvable problems in computer science. Unlike traditional hard-computing algorithms that heavily rely on concrete data and rigid mathematical models, soft

computing leverages a more flexible approach to problem-solving, akin to the human brain's operation, which allows for approximate reasoning and partial truths [2].

2.1. Constituents of Soft Computing

Soft computing encompasses various techniques and methodologies, including FL, EC, ML, and PR. These components work together to handle the uncertainty, ambiguity, and approximation inherent in many real-world problems [2].

2.1.1. Fuzzy Logic

Fuzzy logic was introduced by Lotfi Zadeh in 1965 and provides a computational paradigm that handles uncertainties in data by using levels of truth rather than rigid binary values of 0 and 1 [8]. It has been utilized to create decision systems capable of predicting risks and managing uncertainties in various fields such as engineering and healthcare [2].

2.1.2. Artificial Neural Networks (ANNs)

ANNs are computational models that mimic the structure and functioning of the human brain. These networks consist of interconnected nodes, or perceptrons, that process information using complex mathematical operations.

Through training, ANNs can adjust parameters based on input data, making them highly flexible and capable of handling high-level problems like pattern recognition, predictive modeling, and data analysis [2]. ANNs have found applications in image recognition, natural language processing, and speech recognition, enhancing the versatility and efficiency of soft computing [8].

2.1.3. Evolutionary Computation

Evolutionary computation includes algorithms that mimic natural processes such as evolution and natural selection to find optimal solutions. Techniques like crossover, mutation, and selection drive genetic programming to diversify data and prevent premature conclusions, ensuring the survival of the fittest solutions within a set [8]. These methods are particularly useful in fields such as data mining, robotics, optimization, and engineering [8].

2.2. Applications of Soft Computing

Soft computing techniques are employed across a wide range of industries to solve complex problems efficiently and cost-effectively. For instance, in the medical field, soft computing assists in image analysis, aiding in the recognition and bifurcation of patterns from medical images and X-rays [2]. In the automobile industry, fuzzy logic is used to develop control systems for engine management, automatic transmissions, and antiskid steering [2]. Moreover, the construction industry utilizes Distributed Artificial Intelligence and fuzzy genetic agents to create intelligent buildings that

can adapt to changing environmental conditions [2]. Soft computing also plays a crucial role in handwritten script recognition, allowing for the translation and sorting of multilingual documents [2].

3. Optimization in Soft Computing

Optimization in soft computing involves the application of algorithms and techniques to find the best solution to a given problem, within defined constraints and in the shortest possible time. This process is crucial in various fields such as engineering, management, and artificial intelligence, where complex decision-making and real-life problem-solving are required.

Soft computing provides low-cost, efficient solutions by leveraging algorithms, databases, Fuzzy Sets, and ANNs. These techniques modify uncertainties and indifferences in data, thereby enhancing the decision-making process and minimizing costs [2].

In contrast to traditional hard-computing methods that rely on precise and logical computations, soft computing aims to produce approximate solutions by emulating the human brain's ability to deal with ambiguities and uncertainties. This approach involves several components, including Fuzzy FL, EC, ML, and PR [8].

The book on optimization in engineering highlights the concept of optimization and its relevance to soft computing techniques. It presents various soft computing methods, sample problems, and related software programs, providing a step-by-step guide for obtaining optimal solutions to engineering and management problems [9]. The book emphasizes the broad application of these techniques in solving real-life complex problems using a heuristic approach, covering multidisciplinary areas such as physics, chemistry, biology, and material science [10].

The use of soft computing techniques is particularly beneficial in industries that require the management of vast resources and computations. For instance, soft computing methods are applied in the power system industry to predict voltage instability using ANNs, providing a low-cost solution to ensure system stability [2]. Moreover, these techniques aid in designing optimal maintenance strategies for electrical transmission networks, thereby improving current preventive maintenance programs [11].

4. Application to ELT

The integration of optimization in soft computing techniques within ELT in Iran has been increasingly recognized for its potential to enhance

educational outcomes. The role of English competence in facilitating the use of computers in educational settings is paramount. For instance, interviewees have highlighted that proficiency in English aids significantly in navigating technological challenges, such as understanding instructions during software installation and usage, which are crucial for integrating technology into teaching practices [5]. This underscores the necessity for English language training for teachers of other subjects, positioning English teachers as pioneers in the technological integration process. Moreover, the pedagogical efficacy of incorporating novel technologies and blended learning patterns into ELT has been supported by research. These methods create a more meaningful learning environment, catering to the evolving educational needs driven by technological advancements [12]. Emerging technologies are becoming a staple in students' daily lives, necessitating an adaptive educational framework that embraces these changes.

The application of soft computing in ELT is further illustrated through various studies on assessment and learning processes. For example, LOA has been shown to be effective in enhancing language acquisition. This approach, which focuses on continuous assessment aligned with learning goals, supports the development of language skills in a more contextualized and engaging manner [6].

Additionally, the evaluation of English for Specific Purposes (ESP) courses for Iranian engineering students has revealed positive outcomes in terms of fostering learner autonomy, providing authentic instructional content, and meeting students' real needs [13]. The adoption of soft computing techniques in these courses has demonstrated potential in creating customized and effective learning experiences that resonate with students' professional and academic requirements.

5. Case Studies in Iran

The study of technical English pedagogy specific to the students of engineering in Iran appears to be limited despite the significant amount of research conducted on the effectiveness of tertiary-level ESP courses within the Iranian academic context [13]. Several investigations have explored the effectiveness of ESP courses across various disciplines (e.g., Abbasian & Mahdavi, 2011; Iranmehr, Atai, & Babaii, 2018; Khoshsima & Khosravani, 2014; Mahdavi Zafarghandi et al., 2014; Malmir & Bagheri, 2019; Mashhadi Heidar & Abassy Delvand, 2015; Mazdayasna & Tahririan, 2008; Mostafaei Alaei & Ershadi, 2017; Sherkatolabbasi & Mahdavi Zafarghandi, 2012;

Zand-Moghadam, Meihami, & Ghiasvand, 2018), but few have specifically focused on the needs of engineering students [13].

The limited studies addressing technical English for engineering students (e.g., Danaye-Tous & Haghghi, 2014; Hatam & Shafiei, 2012) often adopt a fragmental approach, evaluating only specific aspects of ESP courses, such as textbooks or language learning skills [13]. This piecemeal approach overlooks the comprehensive evaluation necessary to address the multifaceted needs of engineering students comprehensively.

A detailed examination of various case studies reveals several key findings. For instance, a statistical comparison was drawn between students' and instructors' attitudes toward ESP courses. The data indicated a significant difference in perceptions, particularly in areas such as content authenticity and the satisfaction of target situation needs ($F(1, 848) = 8.215, p < .0125, \eta^2 = .010$; $F(1, 848) = 9.389, p < .0125, \eta^2 = .011$) [13]. The results showcased that students generally expressed dissatisfaction with the authenticity of the speaking, writing, vocabulary, and grammar materials/tasks used in ESP textbooks, with instructors' attitudes being even more negative [13].

Moreover, the analysis of needs satisfaction revealed that both students and instructors agreed on several aspects of the courses, such as the assessment system, task type, teacher expertise, scaffolding type, place of the course in the curriculum, and training hours [13]. However, there was a notable difference in perception regarding instructional materials, with a significant number of instructors disagreeing that courses utilized both prescribed textbooks and up-to-date online resources [13].

A closer look at the capability of the courses to foster learner autonomy showed that most participants believed the courses aimed to promote autonomous ESP learning by introducing extra materials along with the prescribed content [13]. Despite this, the majority of respondents did not confirm the use of other techniques to promote autonomy [13].

6. Benefits and Challenges

6.1. Benefits

The implementation of LOA in English as a Foreign Language (EFL) classes in Iran has shown several significant benefits. One of the most notable advantages is the enhancement of classroom interaction, collaboration, and engagement among students. According to findings derived from MAXQDA analysis, these

factors are frequently cited by educators as key benefits of LOA [6]. Furthermore, LOA facilitates a more relaxed classroom atmosphere, effectively integrating assessment with teaching and learning processes, and enabling a more comprehensive monitoring of learner progress [6]. Teachers have also pointed out that LOA provides a process-oriented assessment which not only improves teaching and learning but also boosts students' motivation and fairness in evaluations [6].

For instance, Teacher 11 emphasized that LOA's integration of assessment, teaching, and learning significantly enhances classroom interaction and self-monitoring among students [6]. Teacher 36 noted that LOA's individualized approach helps in understanding students' needs, strengths, and weaknesses, thereby creating a productive learning environment [6]. Additionally, Teacher 40 observed that LOA reduces the stress associated with traditional testing and motivates students by allowing them to track their own progress [6]. Collectively, these insights underscore the multifaceted benefits of LOA in fostering a dynamic and supportive EFL classroom environment.

6.2. Challenges

Despite its benefits, the implementation of LOA in EFL classes in Iran is not without challenges. Key issues identified include time constraints, large class sizes, and a lack of proper training for teachers on LOA principles and practices [6]. The exam-oriented culture prevalent in many educational contexts further complicates the adoption of LOA, as it often conflicts with the formative and continuous nature of LOA practices [6].

Teacher 23 highlighted that time limits, class size, and insufficient training are significant barriers to effectively applying LOA [6]. Moreover, there is a recognized bias in some LOA practices, such as self and peer assessments, which can undermine the fairness and reliability of evaluations [6]. The findings align with Alsowat (2022), who also identified similar challenges in the context of Saudi Arabia, pointing out the inadequacies in infrastructure and resources required to support LOA [6]. These challenges reflect broader issues in EFL contexts, where traditional assessment approaches still dominate due to their perceived reliability and ease of implementation [6].

7. Future Prospects

The field of optimization in soft computing holds significant potential for transforming ELT in Iran through innovative assessment approaches and

personalized instruction. Future studies are encouraged to explore various aspects to enhance the effectiveness and acceptance of LOA among Iranian EFL teachers. These aspects include investigating the role of teaching experience and educational background in teachers' understanding of LOA, as well as the influence of cultural factors on their perceptions and practices. Cross-cultural studies may offer valuable insights into these dimensions [6].

Additionally, there is a need to examine specific assessment techniques, such as teacher questioning strategies, within the framework of LOA. Such techniques can provide deeper insights into how teachers can effectively integrate LOA into their instructional practices to foster better learning outcomes [6]. Another promising area for future research is the alignment between EFL students' perceptions of LOA and teachers' actual practices, which could lead to a more cohesive and supportive educational environment [6].

In the realm of data analysis and system design, the development of novel data mining algorithms based on fuzzy functional dependencies and linguistic interpretations is a promising avenue. These algorithms support domain experts, including linguists, in making informed decisions during data analysis by expressing useful knowledge from big data in linguistic terms [14]. The application of these algorithms in intelligent tutoring systems can significantly enhance the quality of English education by personalizing instruction based on individual learning patterns, strengths, and weaknesses [7].

Moreover, the integration of deep learning techniques, such as NC-PSOO with GANs, presents new opportunities for improving teaching quality. These advanced computational methods can optimize the teaching and learning process by analyzing various educational variables, ultimately leading to more effective and tailored educational strategies [7].

8. Conclusions

In conclusion, this paper highlights the potential of soft computing solutions in enhancing ELT in Iran by effectively bridging the gap between technology and pedagogy. By integrating adaptive learning systems, natural language processing, and intelligent tutoring systems, educators can create a more personalized and engaging learning environment that caters to diverse student needs. Regarding the novelties of this work, several scientific and practical advancements can be highlighted that may contribute to improving English language teaching in Iran. These include

personalized learning systems, sentiment and interaction analysis, interactive simulations, data analysis for education, and the development of multimedia content. The findings suggest that leveraging these innovative technologies not only improves language acquisition but also fosters learner autonomy and motivation. Ultimately, this approach paves the way for a more effective and dynamic English language education system in Iran, preparing students for a globalized world.

Acknowledgement

The authors would like to thank all the organizations that provided data for this work.

References

- [1] Parvizi, G.R., Tavakoli, M., Amiryousefi, M. and Rezazadeh, M., 2024. Simulating and evaluating individualized cognitive abilities of Iranian EFL learners in orthography acquisition using multi-layer perceptron neural network–gray wolf optimizer computational model. *Education and Information Technologies*, 29(5), pp. 5753-5806. <https://doi.org/10.1007/s10639-023-11825-2>
- [2] Applications of Soft Computing, WisdomPlexus, 2019. <https://wisdomplexus.com/blogs/applications-soft-computing> [Accessed 05 August 2024]
- [3] Mangaroska, K., Sharma, K., Gašević, D. and Giannakos, M., 2020. Multimodal Learning Analytics to Inform Learning Design: Lessons Learned from Computing Education. *Journal of Learning Analytics*, 7(3), pp. 79-97. <https://doi.org/10.18608/jla.2020.73.7>
- [4] Tohidian, I. and Nodooshan, S.G., 2021. Teachers' engagement within educational policies and decisions improves classroom practice: The case of Iranian ELT school teachers. *Improving Schools*, 24(1), pp. 33-46. <https://doi.org/10.1177/1365480220906625>
- [5] Soft computing, WIKIPEDIA, 2024. https://en.wikipedia.org/wiki/Soft_computing [Accessed 05 August 2024]
- [6] Kumar, K., Roy, S. and Davim, J.P., 2019. *Soft computing techniques for engineering optimization*. 1st Edition, CRC press. <https://doi.org/10.1201/9780429053641>
- [7] Kumar, K., Roy, S. and Davim, J.P., 2019. Chapter: Introduction to Optimization and Relevance of Soft Computing towards Optimal Solution. 1st Edition, CRC press. <http://dx.doi.org/10.1201/9780429053641-1>
- [8] Manhal, M., 2007. Fuzzy logic and its application in linguistics. *Journal of the College of Languages (JCL) Mağallaṯ kulliyyaṯ al-luġāt*, (17), pp. 1-32. <https://jcolang.uobaghdad.edu.iq/index.php/JCL/article/view/474>
- [9] Ashtarian, S., 2016. Integration of ICT in English Classrooms in Iran: Affordances and Problems. <https://blog.nus.edu.sg/celcblog/2016/12/01/integration-of-ict-in-english-classrooms-in-iran-affordances-and-problems> [Accessed 05 August 2024]
- [10] Xu, X., Liu, F. and Liang, H., 2023. Role of situational pedagogics in English teaching in the context of soft computing and fuzzy sets. *Journal of Computational Methods in Sciences and Engineering*, 23(5), pp. 2389-2399. <http://dx.doi.org/10.3233/JCM-226818>
- [11] Derakhshan, A. and Ghiasvand, F., 2022. Demystifying Iranian EFL teachers' perceptions and practices of learning-oriented assessment (LOA): Challenges and prospects in focus. *Language Testing in Asia*, 12(1), p. 55. <http://dx.doi.org/10.1186/s40468-022-00204-2>
- [12] Mostafavi, S., Mohseni, A. and Abbasian, G.R., 2021. The pedagogical efficacy of ESP courses for Iranian students of engineering from students' and instructors' perspectives. *Asian-Pacific Journal of Second and Foreign Language Education*, 6, pp. 1-20. <https://sfleducation.springeropen.com/articles/10.1186/s40862-021-00109-2>
- [13] Alikovich Eshbayev, O., Xamidovich Maxmudov, A. and Urokovich Rozikov, R., 2021, December. An overview of a state of the art on developing soft computing-based language education and research systems: a survey of engineering English students in Uzbekistan. In *Proceedings of the 5th International Conference on Future Networks and Distributed Systems*, pp. 447-452. <https://doi.org/10.1145/3508072.3508160>
- [14] Tian, M., 2024. Evaluation model of English Informatization Teaching Quality in Universities Based on Particle Swarm. *Journal of Electrical Systems*, 20(1), p. 139. <https://doi.org/10.52783/jes.672>



Paper Type (Research paper)

A Novel Approach for Intrusion Detection System in IoT Using Correlation-Based Hybrid Feature Selection and Harris Hawk Optimization Algorithm

Yashar Salami¹, Yaser Ebazadeh², Mehdi Hamrang², Nooshin Allahbakhshi³

¹Department of Computer and Information Technology Engineering, Qazvin Branch, Islamic Azad University, Qazvin, Iran

²Department of Computer Engineering, Germe Branch, Islamic Azad University, Germe, Iran

³Department of Computer and Information Technology Engineering, Khoy Branch, Islamic Azad University, Khoy, Iran

Article Info

Article History:

Received: 2024/11/04

Revised: 2024/12/06

Accepted: 2024/12/16

DOI:

Keywords:

Data Mining, Security,
Machine Learning, Anomaly
Detection, Cybersecurity

*Corresponding Author's Email
Address:
yaser_ebazadeh@yahoo.com

Abstract

With the rapid growth of the IoT, the number of devices connected to various networks has significantly increased. These devices generate vast amounts of data and are often deployed in open and unsecured environments, making them vulnerable to cyber-attacks. Therefore, ensuring the security of IoT networks has become a primary concern for researchers. One of the most effective methods for maintaining network security is using IDS. Intrusion detection monitors and analyzes incoming data to detect suspicious activities and potential attacks. Given the resource constraints of IoT devices and the complexity of the networks, improving the accuracy and efficiency of IDS is crucial. The primary goal of this research is to present a novel and optimized IDS for IoT networks. A hybrid feature selection method has been employed to enhance accuracy and reduce computational complexity, combining correlation-based filtering and wrapper methods using the (HHO) algorithm. In this approach, unnecessary features are removed, and essential features for classification are selected. Simulation results indicate that this method has achieved a 96.46% accuracy, outperforming traditional methods such as DT and SVM while improving false positive and false negative rates.

1. Introduction

With the rapid advancement of technology and the integration of the IoT into daily life, IoT has emerged as one of the most significant and transformative innovations of the past decade[1], [2]. By 2030, around 30 billion devices will be connected to the IoT, being utilized in various aspects of life, from smart homes to smart cities[3], [4]. These devices, equipped with the ability to sense their surroundings, collect data, and perform automated actions, have revolutionized our lives[5], [6]. The connection of these devices to cloud technologies has dramatically

enhanced data storage and analysis processes. The potential of IoT to transform our lives is fascinating[7], [8]. However, along with these advantages, significant security challenges have also emerged.

The heterogeneity of IoT devices, the use of various communication protocols, and the large volume of data processed by these devices make them attractive targets for cyber-attacks[6], [9]. Considering these challenges, IoT devices, especially in sensitive environments such as military sectors, manufacturing industries, and smart grids, are

exposed to numerous risks[10], [11]. Cybercriminals and hackers can easily exploit these vulnerabilities to access sensitive data and critical information. This security concern demands attention, as well as the development of novel approaches to counter IoT threats and vulnerabilities[12], [13].

Security studies have shown that traditional methods for combating these threats in IoT devices are not effective due to resource limitations and the complexity of communications[14][15]. Particularly, with cloud and fog computing that leads to vast amounts of data exchange between various devices and cloud servers, traditional security approaches cannot effectively protect these devices[16][17]. Therefore, the need for lighter and more optimized security methods for detecting and preventing cyber-attacks is strongly felt[10].

IDS has been introduced as one of the most effective security tools in this context[18]. These systems analyze and monitor the behavior patterns of IoT devices and networks to detect any abnormal activities or potential cyber-attacks. IDS can identify complex and unknown attacks using machine learning techniques and optimization algorithms[19]. By analyzing data and detecting anomalous patterns in device behavior, IDS alerts system administrators and allows them to take preventive measures before an attack occurs[20].

In environments such as smart homes and smart cities, where IoT devices are widely used, IDS has become one of the key tools for preventing cyber-attacks. The importance of sensitive information in these environments has made high accuracy and efficiency in detecting intrusions and cyber-attacks a necessity[21]. In this regard, using attack modeling and DL techniques can help organizations better understand the nature of attacks and effectively manage security risks[22].

DL techniques, especially given their high computational power and ability to train complex models, are essential in detecting hidden and intricate patterns in IoT data[23]. Utilizing the vast and diverse data generated by IoT devices, these techniques can effectively

identify threats and attacks and enhance network security[24]Overall, transitioning from traditional methods to modern intrusion detection approaches and relying on advanced security techniques is essential to protecting sensitive information and ensuring the security of IoT devices and networks.

Thus, developing machine learning and DL-based security systems, particularly in the IoT domain, can significantly enhance security and reduce these devices' vulnerabilities. The future of security in the IoT world depends on advancements in intelligent security techniques that can keep pace with the growing number of devices and the complexity of their communications while predicting and mitigating cyber threats.

1.1. Contribution

Ultimately, DL techniques play a significant role in detecting attacks and enhancing network security with their high computational power and ability to train complex models. These techniques allow organizations to identify hidden and intricate patterns within IoT data and use them to prevent cyber-attacks. The transition from traditional methods to modern intrusion detection approaches and the adoption of advanced security techniques, particularly in IoT environments, is essential for safeguarding sensitive information and ensuring network security.

1.2. Paper organization

This article's general structure will be described in detail so that the reader can understand the arrangement and contents of each section well. The second part examines and analyzes related works. In this section, the previous research on the subject will be reviewed, and the strengths and weaknesses of each will be carefully examined. This literature review will help the reader to gain a better understanding of the current situation and the need for further research. In the third part, the research proposal is presented. This section describes in detail the methods and techniques chosen to solve the problem in question. Also, the advantages of using these methods compared to the previous methods and the innovations used in this plan are explained in detail. The fourth part deals with simulation and results. Here, the performed simulations are described and their results are analyzed. Also, the graphs and tables used in this section help to show the research findings visually and provide a deeper analysis. Finally, the fifth

section is dedicated to summarizing and concluding. In this section, the key results of the research are reviewed with an emphasis on their importance, and suggestions for future research are also presented. This section, as a summary of the discussed topics, allows the reader to form their own opinions related to the findings and conclusions.

1.3. Symbols

The symbols used in this article are listed in Table 1.

Table 1: symbols used in the article.

<i>IoT</i>	<i>Internet of Things</i>
<i>IDS</i>	<i>Intrusion Detection Systems</i>
<i>HHO</i>	<i>Harris Hawk Optimization</i>
<i>SVM</i>	<i>Support Vector Machine</i>
<i>DT</i>	<i>Decision Tree</i>
<i>PCA</i>	<i>Principal Component Analysis</i>
<i>RBM</i>	<i>Restricted Boltzmann Machines</i>
<i>DL</i>	<i>deep learning</i>
<i>DDOS</i>	<i>Distributed Denial of Service</i>
<i>CFS</i>	<i>Correlation-Based Feature Selection</i>
<i>KNN</i>	<i>K-Nearest Neighbors</i>
<i>TP</i>	<i>True Positive</i>
<i>FP</i>	<i>False Positive</i>
<i>TN</i>	<i>True Negative</i>
<i>FN</i>	<i>False Negative</i>

2. Related work

In [25], a novel IoT intrusion detection approach is presented based on a migration DL model. This method utilizes a feature extraction algorithm to enhance system performance by combining DL and intrusion detection technology. Experiments conducted using the KDD CUP 99 dataset have shown that the proposed algorithm outperforms other algorithms with just 10% of the training data. Empirical results indicate that this approach reduces detection time while providing higher accuracy in identifying attacks. These findings demonstrate that the migration DL model and feature extraction algorithms can significantly improve security in IoT networks and smart cities.

In [26], a framework for intrusion detection based on RBM is proposed. This framework, utilizing RBM—an artificial neural network—can learn high-level features from raw data without supervision. Experimental results from real data collected from a smart water distribution factory showed that this

framework has remarkable performance in detecting attacks and has significantly improved system efficiency. This framework indicates that RBM can be an effective tool for enhancing security in the IoT.

In [27], PCA is employed for feature dimensionality reduction and combined with ensemble-based classifiers to predict intrusion attacks. This research, focusing on intelligent hospitals and medical devices, demonstrated that using PCA and classifiers could enhance the security of smart networks and prevent intrusions in critical systems. The KDDCup'99 dataset was utilized to test this method, yielding significant improvements in attack prediction and enhanced Internet of Medical Things security.

In [28], a group intrusion strategy based on cyber intelligence is introduced, utilizing a combination of machine learning algorithms such as Random Forest, Bayesian Network, C5.0, and CART for detecting botnet attacks in IoT networks. This strategy, employing a cyber intelligence framework and group learning, successfully identified various attacks with high accuracy. Results indicated that applying this framework in smart cities significantly improved detection rates and reduced false favorable rates, contributing to enhanced security in IoT networks.

In [29], methods for improving the performance of IDS against malicious attacks are examined. This method enhances the IDS performance by utilizing adversarial retraining to cope with attacks. Experimental results showed that adversarial retraining could increase the detection accuracy of the IDS to over 99% against malicious attacks, thereby strengthening security in smart cities.

In [30], artificial intelligence techniques for enhancing security in smart cities are explored. Security and privacy issues have gained more attention as the use of information technologies in data management increases in intelligent cities. This study revealed that employing sophisticated artificial intelligence techniques could help improve monitoring, increase security, and protect data in intelligent city networks.

In [31], a method for detecting DDoS attacks in IoT networks is presented. This method

utilizes feature selection to reduce data dimensionality and improve IDS performance. This method selects appropriate features for attack detection using multi-objective optimization and employs extreme learning machines to enhance detection accuracy. Results indicate that this method can effectively increase security in IoT networks. In [32], a DL--based intrusion detection system for IoT devices is introduced. This system employs a four-layer deep network to identify malicious traffic and has demonstrated satisfactory performance in attack detection with an accuracy of 93.74%. This system can function independently of communication protocols and enhance security in IoT networks.

3. Proposed Scheme

This research proposes a combined filter-wrapper approach for feature selection aimed at intrusion detection in IoT networks. This approach consists of two main stages. The first stage is the filter stage, where features related to the class label are identified using correlation measures between the features and the class labels. In this stage, features that correlate with a certain threshold are selected as important and influential features. This stage aims to reduce the dimensionality of the data and select prominent features using statistical criteria. In the next stage, the wrapper stage, the HHO algorithm is used to select a quasi-optimal subset of the selected features. Inspired by the group hunting behavior of Harris hawks, the HHO algorithm finds the best combination of features. This algorithm acts as a population-based optimization method, aiming to improve feature selection through optimization.

The combined filter-wrapper approach in this study enhances the accuracy and speed of intrusion detection in IoT. In this approach, the filter stage helps reduce computational complexity by identifying effective features. Then, the wrapper stage utilizes a machine learning model and the HHO algorithm to refine the selected features. The objective of these two stages is to improve the efficiency and accuracy of the final model simultaneously.

This combined approach's advantage is that it simultaneously employs filter and wrapper methods, reducing data dimensionality and increasing detection accuracy. This combination can be suitable for problems with complex datasets and numerous features, as it simultaneously

reduces computational complexity and enhances model accuracy.

Feature selection algorithms generally select a subset of essential features from large datasets. Their primary goals are to reduce data dimensions, improve the performance of machine learning models, and decrease computational complexity. Eliminating unnecessary or redundant features results in faster model training and increased accuracy. Moreover, appropriate feature selection helps reduce the risk of overfitting and improves the model's generalization on new data.

3.1. Filter-Based Feature Selection

In the feature selection process, there are usually two main stages:

1. **Initial Feature Selection:** In this stage, the goal is to find the smallest subset of features that results in the most minor classification error. Various methods are available, including filter, wrapper, and hybrid methods.
2. **Model Evaluation:** In this stage, the model is trained using the selected features, and its performance is evaluated using test data. This stage aims to reduce the training samples' classification error and improve the test samples' prediction accuracy.

Selecting essential and useful features can significantly improve the performance and accuracy of machine learning models. Therefore, feature selection is considered a fundamental step in the machine learning process.

Feature selection methods fundamentally improve the performance of machine learning models by removing unnecessary or redundant features, reducing data dimensions, or increasing the models' generalization capability. These methods help models better identify important patterns in the data and prevent overfitting to the training data.

In supervised feature selection, a machine learning model is trained using labeled data, and important features for accurately predicting labels are determined. These methods are typically based on criteria such as information gain, feature importance, or prediction error.

In unsupervised feature selection, data is used without labels, and only essential and separable features are validated. These methods are generally performed based on dimensional analysis or pattern recognition capabilities.

In semi-supervised feature selection, labeled and unlabeled data are used to select important features for class separation. These methods are often suitable

for addressing issues like the scarcity of labeled samples.

3.2. Feature Selection Based on Correlation

The CFS algorithm is one of the most powerful and popular methods for feature selection in classification problems. Due to its simplicity and high efficiency, this algorithm is widely used in many machine-learning issues. The main objective of CFS is to identify a subset of features that not only have a strong correlation with target classes but also exhibit a low correlation with one another. This reduces redundant information among features and helps selected features better differentiate between classes.

The algorithm initially ranks all features based on their correlation with the classes. After that, a subset of features with a high correlation to the classes and a low correlation with each other is selected. These features are recognized as a set of features with "best generalization" and play a key role in enhancing the performance of classification models. These features make the models more accurate and increase their ability to detect and classify data.

The advantage of using the CFS algorithm as a filter for feature selection lies in its ability to identify and eliminate irrelevant features or those with a high correlation to other features. This process reduces model complexity and increases efficiency, as only features related to the classes remain independent.

The equation (1) used in this algorithm to evaluate the quality of a subset of features is as follows:

$$Merit_s = \frac{kr_{cf}}{\sqrt{k+(k+1)r_{ff}}} \quad (1)$$

Where Merit is the exploratory "merit" of a feature subset (S) containing (k) features, (rcf) is the average class correlation of the feature where (rff) is the average feature-feature correlation. The numerator of this equation acts as a class prediction index using the selected features, and the denominator indicates redundancy among the features.

CFS employs Pearson correlation to evaluate the features. This method standardizes features and ranks them based on their correlation with target classes. Features that have little correlation with the classes typically provide useless or inefficient information for classification and are, therefore, eliminated. Additionally, redundant features that correlate highly with one or more other features are discarded from the model to enhance its efficiency and accuracy. This process increases accuracy and makes classification models more efficient in processing data and correctly predicting classes.

3.3 Wrapper-Based Feature Selection

Wrapper-based methods for feature selection utilize classification models to evaluate and select a subset of features that yield the best performance for the model. These wrapper methods employ a search strategy to find the optimal subset of features. Such strategies may include forward sequential search, backward sequential search, genetic algorithms, etc. A subset of features is selected at each step of the search, and a learning model is trained on this subset. The model's performance is assessed based on accuracy, sensitivity, specificity, etc. The subset of features that provides the best model performance is then selected as the final subset.

These methods directly leverage the learning model for evaluation, enabling them to choose features that contribute optimally to the model's performance. The model becomes more straightforward and interpretable by eliminating unnecessary features and retaining important ones. Selecting fewer features reduces the time and resources required for training and using the model. However, these methods require repeated execution of learning algorithms on various subsets of features, which can be very time-consuming and costly due to the significant resources needed for searching and evaluating.

Wrapper-based methods can improve the performance of machine learning models. Still, they require a thorough and meticulous evaluation of feature subsets and the selection of appropriate modeling and assessment techniques. These methods can help reduce model complexity, improve interpretability, and enhance prediction accuracy.

3.4 Feature Selection Based on HHO

This research employs the HHO algorithm to select optimal features that correlate highly with class labels. The primary objective of this approach is to identify significant features from datasets related to the IoT. The HHO algorithm aims to optimize the feature selection process by considering two fundamental goals: reducing the number of features and enhancing the accuracy of intrusion detection. The process begins with the algorithm analyzing the existing features, attempting to identify key and important features using heuristic methods. Subsequently, through computationally implemented learning and optimization techniques, the algorithm moves towards finding a subset of features that exhibits the most minor complexity while maintaining the highest accuracy in intrusion detection.

Given that the feature selection problem is classified as NP-Hard, the HHO utilizes specific optimization methods to find a genuinely optimal solution. This method is inspired by the natural behaviors of Harris hawks, which employ a combination of random and guided search strategies to discover the best features.

3.5 Fitness Function

The fitness function plays a critical role in evolutionary and metaheuristic algorithms, serving as a key stage for assessing the quality and performance of each solution within a population. This function evaluates the performance of each solution and determines its success in addressing the problem at hand. In other words, the value of the fitness function is determined based on the selected features from the datasets, indicating how closely each solution aligns with the desired objective.

This function automatically computes the target value based on criteria such as the prediction error of other defined metrics. Solutions yielding a higher fitness value are recognized as superior candidates and are selected for subsequent iterations of the evolutionary algorithm. Ultimately, solutions that achieve the best fitness values are designated as the final and optimal results and utilized to solve the problem effectively. This iterative process contributes to creating an evolutionary cycle, facilitating continuous improvement of solutions and optimizing the resolution of complex issues.

In the next phase, each solution is evaluated based on defined objectives and the fitness function. This assessment employs a proposed fitness function that combines the feature selection rate with the prediction error of intrusion in the IoT context. The fitness function shows question 2 is defined as follows:

$$\text{Minimize } F(x) = \begin{cases} f_1(x) = \frac{L}{A}, & L \in A, A \in \mathbb{R}^+ \\ f_2(x) = \frac{FP+FN}{P+N}, & (P+N) \in \mathbb{R}^+ \end{cases} \quad (2)$$

In this equation (2), A represents the total number of features, while L indicates the number of selected features. These

calculations intertwine with issues related to security and safety within systems. The variables TP, FP, TN, and FN denote the actual number of healthy nodes, nodes erroneously identified as intrusions, nodes correctly identified as intrusions, and nodes mistakenly recognized as healthy, respectively.

The set P encompasses the sum of TP and TN, while the set N consists of the sum of FP and FN. These variables are utilized to evaluate the accuracy and quality of a model or security system, ensuring that the feature selection process effectively enhances the model's performance while minimizing the risk of misclassification.

4 Simulation

This paper presents an innovative IoT intrusion detection and prediction method. This method is based on a combination of feature selection using a filter-wrapper approach based on correlation and using HHO to classify the data. This section introduces the standard dataset used for intrusion detection in IoT. This dataset, named BotNeTIoT, is extracted from the UCI database and contains information related to over 50,000 nodes. Some nodes are classified as healthy, while others are classified as intrusive.

The main objective of using this dataset is to compare the proposed method for intrusion detection in IoT with other existing methods. This comparison aims to evaluate the improvements the proposed method can bring in terms of the accuracy and performance of intrusion detection models. During the implementation of this method, challenges such as data imbalance and the high volume of existing features in the dataset arise. Therefore, selecting key and effective features can significantly improve the accuracy and efficiency of classification algorithms and machine learning, assisting in more precise intrusion detection in IoT. Table 2 shows details of the features available in this dataset.

Table 2: Features the dataset.

Feature	Feature Number	Feature Title	Feature Number
HH_L0.1_pcc	13	MI_dir_L0.1_weight	1
HH_jit_L0.1_weight	14	MI_dir_L0.1_mean	2
HH_jit_L0.1_mean	15	MI_dir_L0.1_variance	3
HH_jit_L0.1_variance	16	H_L0.1_weight	4
HpHp_L0.1_weight	17	H_L0.1_mean	5
HpHp_L0.1_mean	18	H_L0.1_variance	6
HpHp_L0.1_std	19	HH_L0.1_weight	7
HpHp_L0.1_magnitude	20	HH_L0.1_mean	8
HpHp_L0.1_radius	21	HH_L0.1_std	9
HpHp_L0.1_covariance	22	HH_L0.1_magnitude	10
HpHp_L0.1_pcc	23	HH_L0.1_radius	11
Label	24	HH_L0.1_covariance	12

4.1. Dimensionality Reduction Based on High Correlation

In this research, a dataset available in the UCI repository that contains 24 different features is utilized. These features are listed in Table 2. Since not all features may be equally effective in intrusion detection and using all of them could increase model complexity, selecting a subset of useful features is essential. In this study, a CFS approach is employed to optimize the selection of features.

Pearson correlation is a statistical method that examines the linear relationship between two variables. This correlation can take values ranging from -1 to 1, where values close to 1 indicate a direct and positive linear relationship, values close to -1 indicate a reverse relationship and values close to 0 indicate no linear relationship between the variables. In this study, Pearson correlation analysis is used between the features and class labels to identify and select features that have the highest correlation with the class label. This approach improves intrusion detection accuracy and helps reduce data dimensions.

Table 2 presents the results of the Pearson correlation analysis. These results illustrate how selecting more relevant features can enhance the performance of classification models and intrusion prediction in IoT

networks. Table 3 discusses the importance of the correlation between features and class labels in an analytical model. The correlation between features and class labels falls from -1 to 1. A positive correlation indicates a direct relationship between a feature and the class label; as the feature's value increases, the class label's value also increases. Conversely, a negative correlation indicates a reverse relationship, where an increase in the feature value leads to a decrease in the class label value.

Bar charts like Figure 1 are utilized to visualize these relationships better. Based on the calculated correlation, these charts graphically display the relationship between features and the class label. These charts help researchers identify important features that correlate significantly with the class label. By selecting these key features significantly related to the class label, the analytical model can provide better accuracy and performance, ultimately enhancing the model's effectiveness.

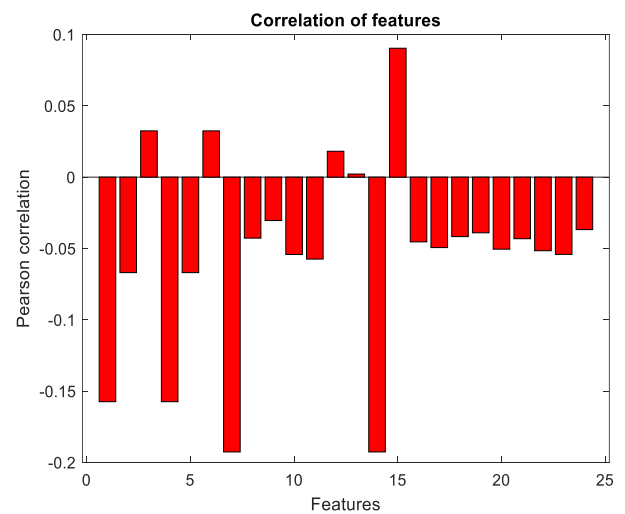


Fig 1. Correlation between features and class label.

In Figure 1, various features exhibit either a positive or negative correlation with the class label and the degree of this correlation varies across different features. This diversity in correlation can play a crucial role in selecting the most suitable features for analytical models. For features with positive correlation, the closer the correlation value is to 1, the stronger the relationship between the feature and the class label. This indicates that if a feature's correlation with the class label

approaches 1, it is considered an effective and valuable element for model inclusion.

Conversely, features with negative correlation are also significant. The closer a negative correlation value is to -1, the stronger the inverse relationship between the feature and the class label. In other words, if the negative correlation of a feature with the class label approaches -1, that feature may also improve the model's performance. The proposed method establishes a specific threshold to identify valuable features. Features that have a positive correlation and whose correlation values exceed the average positive correlation of other features are selected as valuable. Similarly, features with negative correlation and whose correlation values fall below the average negative correlation of other features are also considered useful. Selecting these key features can enhance the accuracy and efficiency of statistical and predictive models. Table 4 presents the valuable features and their correlation values with the class label. The selection of 12 features with the highest correlation with the class labels from the 24 available features in the dataset can help improve the performance and efficiency of the statistical and predictive models. These selected high-correlation features will act as useful and impactful inputs for the models, and by eliminating the less important features, the accuracy and speed of the models will increase. This careful feature selection can lead to an overall enhancement in the performance of the prediction and decision-making systems.

As shown in Table 4, the features that meet the defined threshold criteria are identified as key and valuable features and selected for use in statistical and predictive models. This intelligent and precise feature selection significantly improves model accuracy and efficiency, enhancing predictive systems' performance. Additionally,

The show Figure 2 depicts a **correlation matrix**, which is one of the key tools in data analysis and is used to examine the relationships among variables in a dataset. This matrix, displayed as a square table, provides the Pearson correlation coefficients for each pair of variables. These coefficients help researchers understand how changes in one variable might relate to changes in another.

The matrix is presented both numerically and visually, allowing for a quick and intuitive understanding of patterns and relationships between variables.

The correlation coefficients, which range from [-1, +1], describe the **strength and direction of the linear relationship** between two variables:

- **Positive values (close to +1):** Indicate a strong direct correlation, meaning that as the value of one variable increases, the value of the other variable also increases proportionally. For example, if two economic variables such as "income" and "consumer spending" have a correlation close to +1, an increase in income is likely to result in an increase in spending.
- **Negative values (close to -1):** Indicate a strong inverse correlation, meaning that as the value of one variable increases, the value of the other variable decreases. For instance, in environmental studies, there might be a strong negative correlation between "rainfall levels" and "drought percentage."
- **Values close to zero:** Suggest no linear relationship or a very weak relationship between two variables. This means that changes in one variable have no significant impact on the other. Such relationships might be random or influenced by unrelated factors.



Fig 2. Correlation of selected features.

4.2 Feature Subset Selection Based on HHO

The HHO algorithm is primarily used to find optimal solutions across various problems. As an advanced search method, it explores different solutions to provide an optimal combination of features for classification and prediction tasks. This study evaluates the HHO algorithm based on classification error, which serves as a criterion for assessing solutions. Utilizing this algorithm, an optimal feature set is identified for classifying and predicting intrusions in the IoT. The HHO offers the most efficient combination of intrusion detection features and IoT prediction based on the available characteristics. This method can significantly enhance the performance and efficiency of predictive and classification models while determining the best feature combinations for various issues. Table 5 illustrates the optimal solution selected by the Harris Hawk.

Table 5: Optimal Solutions in HHO

<i>Feature Number</i>	<i>Feature Name</i>
2	<i>H_L0.1_weight</i>
4	<i>H_L0.1_mean</i>
5	<i>HH_jit_L0.1_mean</i>
10	<i>HH_L0.1_magnitude</i>
15	<i>MI_dir_L0.1_mean</i>

According to **Table 5**, the HHO significantly improved intrusion detection within IoT by selecting seven essential features for use in classifying training samples and predicting test samples. The algorithm evaluates its performance using classification error as the assessment criterion, calculated based on a fitness function. The convergence chart presented in **Figure 3** illustrates the reduction in the fitness function and the classification error for intrusions achieved through the solutions provided by the HHO algorithm. As the progresses, both the fitness function value and classification error improve, indicating the effectiveness and efficient performance of this algorithm in detecting intrusions within the IoT environment employing the optimal features selected by the HHO and accurately calculating the classification error as the primary evaluation criterion can enhance classification and intrusion prediction

performance in IoT, facilitating more accurate and reliable predictions.

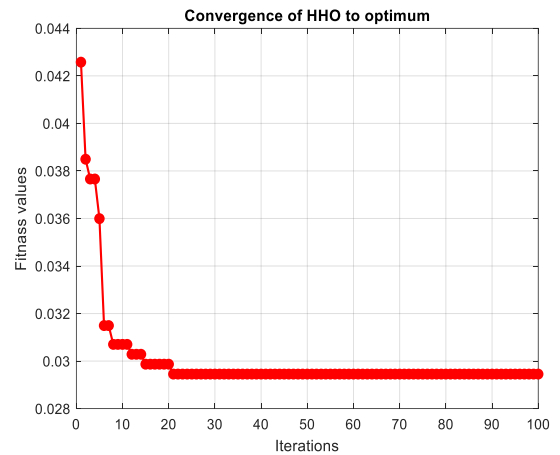


Fig 3. Convergence of the fit function.

As depicted in **Figure 3**, the HHO effectively solves feature subset selection issues for intrusion detection in the IoT. With increased iterations and improved performance, the algorithm approaches the optimal fitness function, yielding high accuracy. For instance, after 100 iterations, the algorithm achieves a fitness function value of **0.029**, demonstrating improved performance and high accuracy. These results indicate that the HHO reaches the best possible solution for the feature subset selection problem and minimizes the classification error to the lowest feasible extent. The convergence chart also shows that the algorithm moves towards the optimal point, reducing the classification error, which reflects better performance and higher accuracy in detecting intrusions within the IoT. This algorithm can substantially improve intrusion detection and prediction, enhancing preventive and security measures. This algorithm leverages the combined strength of the HHO technique and the computational power of s to achieve optimal and efficient solutions for complex issues. This makes the HHO an effective and efficient tool for detecting and preventing intrusions in the IoT.

4.3 Evaluation of the Proposed Method

Various assessment metrics are utilized to evaluate the proposed method for analyzing healthy and intrusion nodes. These metrics are derived from the confrontation between the

predicted class labels for healthy and intrusion nodes and their actual class labels, as delineated in the confusion matrix. The primary evaluation metrics applicable for binary classification methods include:

1. **Accuracy:** Equation (3) represents the ratio of correctly classified nodes to the total number of nodes[33].

$$Accuracy = \frac{TP + TN}{TP+TN+FP+FN} \quad (3)$$

2. **Sensitivity or Recall:** Equation (4) indicates the proportion of actual positive nodes that are correctly classified[34].

$$Recall = \frac{TP}{TP+FN} \quad (4)$$

3. **Specificity:** Equation (5) metric shows the proportion of actual negative nodes that are accurately classified[35].

$$Precision = \frac{TP}{TP+FP} \quad (5)$$

F1-Score: Equation (6) The F1-score is the geometric mean of precision and sensitivity, serving as a combined measure of the model's accuracy and information content[36].

$$F - measure = \frac{2*Precision*Recall}{Precision+Recall} \quad (6)$$

These metrics assist in evaluating the classification method's performance on the test data, allowing us to quantify its accuracy and effectiveness. The confusion matrix, along with the parameters TP, FP, TN and FN, provides a foundational framework for calculating these metrics and comprehensively understanding the classification method's performance. Figure 4 show the pseudocode of the proposed method.

According to Table 6, the proposed method, which integrates features based on Pearson correlation and Harris Hawk optimization, achieves commendable results for the evaluation metrics during testing with new data and nodes. These results indicate that the method performs well in detecting and classifying the targeted cases.

Furthermore, the KNN classification method demonstrates improved performance compared to the other examined methods. This suggests that KNN can be a robust and effective data classification and pattern recognition approach, yielding better results than alternative methods.

4.4 Comparison of the Proposed

In this section, after assessing the performance of the proposed method against the test datasets, the results obtained from this method—which comprises a combination of CFS and the Harris Hawks alongside classification methods—are compared with previous approaches under identical conditions on a specific dataset concerning intrusion detection in IoT. The significance of this comparison stems from the challenges associated with heart disease detection, tracking, and identifying heart patients amidst heterogeneous datasets, primarily due to issues like imbalanced class distributions. Consequently, accuracy metrics that illustrate the relationship between actual samples and those identified by machine learning models serve as the optimal criterion for evaluating the proposed method. Figure 5 presents the bar chart corresponding to the accuracy metric in the proposed method juxtaposed with various classification methods and prior approaches. This method enhances the performance and accuracy of data classification by selecting essential and practical features based on their correlations. Utilizing the Harris Hawks empowers this approach to optimize and refine the performance of the KNN classification model. The selection of significant features through correlation analysis aids in improving the accuracy and precision of data classification, resulting in noteworthy improvements compared to previous methods. The results yielded by this proposed method with prior classification techniques indicate that CFS and the Harris Hawks facilitate superior performance and evident enhancements in the accuracy and effectiveness of data classification. These comparisons demonstrate that the proposed method can deliver more optimal and accurate data classification results while improving classification techniques' performance.

5. Conclusion

The IoT encompasses systems composed of various devices connected to the internet, facilitating information exchange between them. These devices can extensively and adaptively change based on user needs. Security and privacy challenges have markedly increased with the rising utilization of IoT devices. Connected IoT devices access sensitive information, thereby heightening the likelihood of cyberattacks. Thus, developing IDS to secure these devices and their communications is critical. IDS in IoT are designed to identify and prevent intrusions and attacks targeting devices and networks. Typically, these systems utilize various analytical algorithms and models to monitor activities at either local or cloud levels. They play a vital role in enhancing security and protecting devices and data linked to the IoT, assisting organizations in mitigating security threats and preventing attacks. Given the substantial volume of data and the diverse features of nodes within IoT networks, employing feature selection methods to increase the accuracy of IDS is essential. In this research, an intrusion detection system based on a combined filter-wrapper feature selection method, utilizing correlation analysis and the Harris Hawks, has been proposed. Results indicate that this method has optimized the performance of the KNN classification model, achieving an accuracy of 96.46%. The selection of significant features based on their correlations has improved accuracy in data classification, yielding substantial results compared to previous approaches. The comparison of this method with other classification techniques has revealed that the integration of CFS and Harris Hawks optimization results in significant improvements in the accuracy and performance of data classification systems, showcasing the ability to deliver more precise and optimized results in data classification.

References

- [1] M. Wazid, P. Bagga, A. K. Das, S. Shetty, J. J. P. C. Rodrigues, and Y. Park, "AKM-IoV: Authenticated Key Management Protocol in Fog Computing-Based Internet of Vehicles Deployment," *IEEE Internet Things J.*, vol. 6, no. 5, pp. 8804–8817, 2019.
- [2] Y. Salami, V. Khajehvand, and E. Zeinali, "SOS-FCI: a secure offloading scheme in fog–cloud-based IoT," *J. Supercomput.*, vol. 80, no. 1, pp. 570–600, 2024, doi: 10.1007/s11227-023-05499-3.
- [3] X. Mu and M. F. Antwi-Afari, "The applications of Internet of Things (IoT) in industrial management: a science mapping review," *Int. J. Prod. Res.*, vol. 62, no. 5, pp. 1928–1952, 2024.
- [4] M. Sugar and I. H. Berkovitz, "Fog Computing Conceptual Model," *Adolesc. Psychiatry (Hilversum)*, vol. 1, no. 2, pp. 169–178, 2011, doi: 10.2174/2210677411101020169.
- [5] M. Nassereddine and A. Khang, "Applications of Internet of Things (IoT) in smart cities," in *Advanced IoT technologies and applications in the industry 4.0 digital economy*, CRC Press, 2024, pp. 109–136.
- [6] Y. Salami, V. Khajehvand, and E. Zeinali, "SAIFC: A Secure Authentication Scheme for IOV Based on Fog-Cloud Federation," *Secur. Commun. Networks*, vol. 1, pp. 1–19, 2023.
- [7] A. Rajagopalan *et al.*, "Empowering power distribution: Unleashing the synergy of IoT and cloud computing for sustainable and efficient energy systems," *Results Eng.*, p. 101949, 2024.
- [8] Y. Salami, Y. Ebazadeh, and V. Khajehvand, "CE-SKE: cost-effective secure key exchange scheme in Fog Federation," *Iran J. Comput. Sci.*, vol. 4, no. 3, pp. 1–13, 2021.
- [9] A. Souri, M. Norouzi, and Y. Alsenani, "A new cloud-based cyber-attack detection architecture for hyper-automation process in industrial internet of things," *Cluster Comput.*, vol. 27, no. 3, pp. 3639–3655, 2024.
- [10] Y. Salami and S. Hosseini, "BSAMS: Blockchain-Based Secure Authentication Scheme in Meteorological Systems," *Nivar*, vol. 47, no. 120–121, pp. 181–197, 2023.
- [11] Y. Salami, F. Taherkhani, Y. Ebazadeh, M. Nemati, V. Khajehvand, and E. Zeinali, "Blockchain-Based Internet of Vehicles in Green Smart City: Applications and Challenges and Solutions," *Anthropog. Pollut.*, vol. 7, no. 1, pp. 87–96, 2023.
- [12] S. C. Vetrivel, R. Maheswari, and T. P. Saravanan, "Industrial IOT: Security Threats and Counter Measures," in *Communication Technologies and Security Challenges in IoT: Present and Future*, Springer, 2024, pp. 403–425.
- [13] Y. Salami, V. Khajehvand, and E. Zeinali, "Cryptographic Algorithms: A Review of the Literature, Weaknesses and Open Challenges," *J. Comput. Robot.*, vol. 16, no. 2, pp. 46–56, 2023.
- [14] U. D. Maiwada, S. A. Imran, K. U. Danyaro, A. A. Janisar, A. Salameh, and A. B. Sarlan, "Security Concerns of IoT Against DDoS in 5G Systems," *Int. J. Electr. Eng. Comput. Sci.*, vol. 6, pp. 98–105, 2024.
- [15] Y. Salami, V. Khajehvand, and E. Zeinali, "A new secure offloading approach for internet of vehicles in fog-cloud federation," *Sci. Rep.*, vol. 14, no. 1, p. 5576, 2024.
- [16] Y. Salami and V. Khajehvand, "SMAK-IOV: Secure Mutual Authentication Scheme and Key Exchange Protocol in Fog Based IoV," *J. Comput. Robot.*, vol. 13, no. 1, pp. 11–20, 2020.
- [17] Y. Salami, V. Khajehvand, and E. Zeinali, "LSMAK-IOV: Lightweight Secure Mutual AKE Scheme in

- Fog-Based IoV,” in *2024 10th International Conference on Artificial Intelligence and Robotics (QICAR)*, IEEE, 2024, pp. 1–5.
- [18] Z. Wang, J. Li, S. Yang, X. Luo, D. Li, and S. Mahmoodi, “A lightweight IoT intrusion detection model based on improved BERT-of-Theseus,” *Expert Syst. Appl.*, vol. 238, p. 122045, 2024.
- [19] O. B. J. Rabie, S. Selvarajan, T. Hasanin, A. M. Alshareef, C. K. Yogesh, and M. Uddin, “A novel IoT intrusion detection framework using Decisive Red Fox optimization and descriptive back propagated radial basis function models,” *Sci. Rep.*, vol. 14, no. 1, p. 386, 2024.
- [20] E. Altulaihian, M. A. Almaiah, and A. Aljughaiman, “Anomaly Detection IDS for Detecting DoS Attacks in IoT Networks Based on Machine Learning Algorithms,” *Sensors*, vol. 24, no. 2, p. 713, 2024.
- [21] M. M. Inuwa and R. Das, “A comparative analysis of various machine learning methods for anomaly detection in cyber attacks on IoT networks,” *Internet of Things*, vol. 26, p. 101162, 2024.
- [22] A. Aldhaheri, F. Alwahedi, M. A. Ferrag, and A. Battah, “Deep learning for cyber threat detection in IoT networks: A review,” *Internet Things cyber-physical Syst.*, vol. 4, pp. 110–128, 2024.
- [23] N. O. Aljehane et al., “Golden jackal optimization algorithm with deep learning assisted intrusion detection system for network security,” *Alexandria Eng. J.*, vol. 86, pp. 415–424, 2024.
- [24] C. Hazman, A. Guezzaz, S. Benkirane, and M. Azrou, “Enhanced ids with deep learning for iot-based smart cities security,” *Tsinghua Sci. Technol.*, vol. 29, no. 4, pp. 929–947, 2024.
- [25] D. Li, L. Deng, M. Lee, and H. Wang, “IoT data feature extraction and intrusion detection system for smart cities based on deep migration learning,” *Int. J. Inf. Manage.*, vol. 49, pp. 533–545, 2019.
- [26] A. Elsaedy, K. S. Munasinghe, D. Sharma, and A. Jamalipour, “Intrusion detection in smart cities using Restricted Boltzmann Machines,” *J. Netw. Comput. Appl.*, vol. 135, pp. 76–83, 2019.
- [27] T. Saba, “Intrusion detection in smart city hospitals using ensemble classifiers,” in *2020 13th International Conference on Developments in eSystems Engineering (DeSE)*, IEEE, 2020, pp. 418–422.
- [28] E. M. Onyema, S. Dalal, C. A. T. Romero, B. Seth, P. Young, and M. A. Wajid, “Design of intrusion detection system based on cyborg intelligence for security of cloud network traffic of smart cities,” *J. Cloud Comput.*, vol. 11, no. 1, p. 26, 2022.
- [29] M. M. Rashid et al., “Adversarial training for deep learning-based cyberattack detection in IoT-based smart city applications,” *Comput. Secur.*, vol. 120, p. 102783, 2022.
- [30] M. Abdedaime, A. Qafas, M. Jerry, and A. Guezzaz, “A KNN-based intrusion detection model for smart cities security,” in *International Conference on Innovative Computing and Communications: Proceedings of ICICC 2022, Volume 3*, Springer, 2022, pp. 265–272.
- [31] M. Roopak, G. Y. Tian, and J. Chambers, “Multi-objective-based feature selection for DDoS attack detection in IoT networks,” *IET Networks*, vol. 9, no. 3, pp. 120–127, 2020.
- [32] A. Awajan, “A novel deep learning-based intrusion detection system for IOT networks,” *Computers*, vol. 12, no. 2, p. 34, 2023.
- [33] J. Li, M. Gao, and R. D’Agostino, “Evaluating classification accuracy for modern learning approaches,” *Stat. Med.*, vol. 38, no. 13, pp. 2477–2503, 2019.
- [34] J. Miao and W. Zhu, “Precision–recall curve (PRC) classification trees,” *Evol. Intell.*, vol. 15, no. 3, pp. 1545–1569, 2022.
- [35] Ž. Vujović, “Classification model evaluation metrics,” *Int. J. Adv. Comput. Sci. Appl.*, vol. 12, no. 6, pp. 599–606, 2021.
- [36] D. Chicco and G. Jurman, “The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation,” *BMC Genomics*, vol. 21, pp. 1–13, 2020.

Table 3: Correlation Values between Features and Class Labels

<i>Correlation</i>	<i>Feature Name</i>	<i>Feature Number</i>	<i>Correlation</i>	<i>Feature Name</i>	<i>Feature Number</i>
0.0022	HH_L0.1_pcc	13	-0.1574	MI_dir_L0.1_weight	1
-0.1927	HH_jit_L0.1_weight	14	-0.0670	MI_dir_L0.1_mean	2
0.0904	HH_jit_L0.1_mean	15	0.0324	MI_dir_L0.1_variance	3
-0.0454	HH_jit_L0.1_variance	16	-0.1574	H_L0.1_weight	4
-0.0494	HpHp_L0.1_weight	17	-0.0670	H_L0.1_mean	5
-0.0417	HpHp_L0.1_mean	18	0.0324	H_L0.1_variance	6
-0.0390	HpHp_L0.1_std	19	-0.1927	HH_L0.1_weight	7
-0.0505	HpHp_L0.1_magnitude	20	-0.0428	HH_L0.1_mean	8
-0.0432	HpHp_L0.1_radius	21	-0.0304	HH_L0.1_std	9
-0.0516	HpHp_L0.1_covariance	22	-0.0542	HH_L0.1_magnitude	10
-0.0541	HpHp_L0.1_pcc	23	-0.0575	HH_L0.1_radius	11
-0.0368	TnBPSrcIP	24	0.0182	HH_L0.1_covariance	12

Table 4: Selected Features.

<i>F-number</i>	<i>F_Name</i>	<i>Correlation</i>	<i>F_number</i>	<i>F_Name</i>	<i>Correlation</i>
1	MI_dir_L0.1_weight	-0.1574	11	HH_L0.1_radius	-0.0575
2	MI_dir_L0.1_mean	-0.0670	14	HH_jit_L0.1_weight	-0.1927
4	H_L0.1_weight	-0.1574	15	HH_jit_L0.1_mean	0.0904
5	H_L0.1_mean	-0.0670	20	HpHp_L0.1_magnitude	-0.0505
7	HH_L0.1_weight	-0.1927	22	HpHp_L0.1_covariance	-0.0516
10	HH_L0.1_magnitude	-0.0542	23	HpHp_L0.1_pcc	-0.0541

Table 6. Average values of the evaluation metrics for various classification methods

<i>Classification Method</i>	<i>F1-Score</i>	<i>Accuracy</i>	<i>Sensitivity</i>	<i>Specificity</i>
KNN	98.27	97.29	99.25	96.46
Neural Networks	63.57	53.47	91.48	86.18
Decision Tree	97.14	97.22	97.07	95.60
Naive Bayes	93.60	95.98	91.48	91.93
Support Vector Machine	85.30	94.32	77.93	79.33


```

#define MAX_FEATURES Numbers
// Step 1: Feature Extraction
// This function takes sensor data and extracts relevant features.
void extract_features(float data[MAX_FEATURES][MAX_FEATURES], int n) {
    // Collect data and extract features
}
// Step 2: Calculate Correlation
// This function computes the correlation matrix based on the extracted features.
void calculate_correlation(float data[MAX_FEATURES][MAX_FEATURES], int n, float
corr[MAX_FEATURES][MAX_FEATURES]) {
    // Calculate the correlation matrix
}
// Step 3: Identify Causal Relationships
// This function identifies causal relationships among the features using the correlation matrix.
void identify_causal_relationships(float corr[MAX_FEATURES][MAX_FEATURES], int n, int
causal[MAX_FEATURES][MAX_FEATURES]) {
    // Identify causal relationships
}

// Step 4: Analyze Mutual Influence
// This function analyzes the mutual influence of features upon each other.
void analyze_mutual_influence(int causal[MAX_FEATURES][MAX_FEATURES], int n, float
influence[MAX_FEATURES][MAX_FEATURES]) {
    // Analyze mutual influences
}

// Step 5: Optimize Causal Model
// This function optimizes the causal model based on the mutual influences.
void optimize_causal_model(int causal[MAX_FEATURES][MAX_FEATURES], int n, float
influence[MAX_FEATURES][MAX_FEATURES], int optimized[MAX_FEATURES][MAX_FEATURES]) {
    // Optimize the causal model
}

// Step 6: Evaluate Model
// This function evaluates the optimized model and compares the results with other methods.
void evaluate_model(int optimized[MAX_FEATURES][MAX_FEATURES], int n) {
    // Evaluate and compare results of the optimized model
}
int main() {
    float sensor_data[MAX_FEATURES][MAX_FEATURES]; // Sensor data
    float correlation_matrix[MAX_FEATURES][MAX_FEATURES]; // Correlation matrix
    int causal_relationships[MAX_FEATURES][MAX_FEATURES]; // Causal relationships
    float mutual_influence[MAX_FEATURES][MAX_FEATURES]; // Mutual influences
    int optimized_causal_relationships[MAX_FEATURES][MAX_FEATURES]; // Optimized causal
relationships
    int num_features =; // Number of features

    // Execute research steps
    extract_features(sensor_data, num_features);
    calculate_correlation(sensor_data, num_features, correlation_matrix);
    identify_causal_relationships(correlation_matrix, num_features, causal_relationships);
    analyze_mutual_influence(causal_relationships, num_features, mutual_influence);
    optimize_causal_model(causal_relationships, num_features, mutual_influence,
optimized_causal_relationships);
    evaluate_model(optimized_causal_relationships, num_features); // Evaluate the model
    return 0;
}

```

Fig 4. pseudocode of the proposed method.

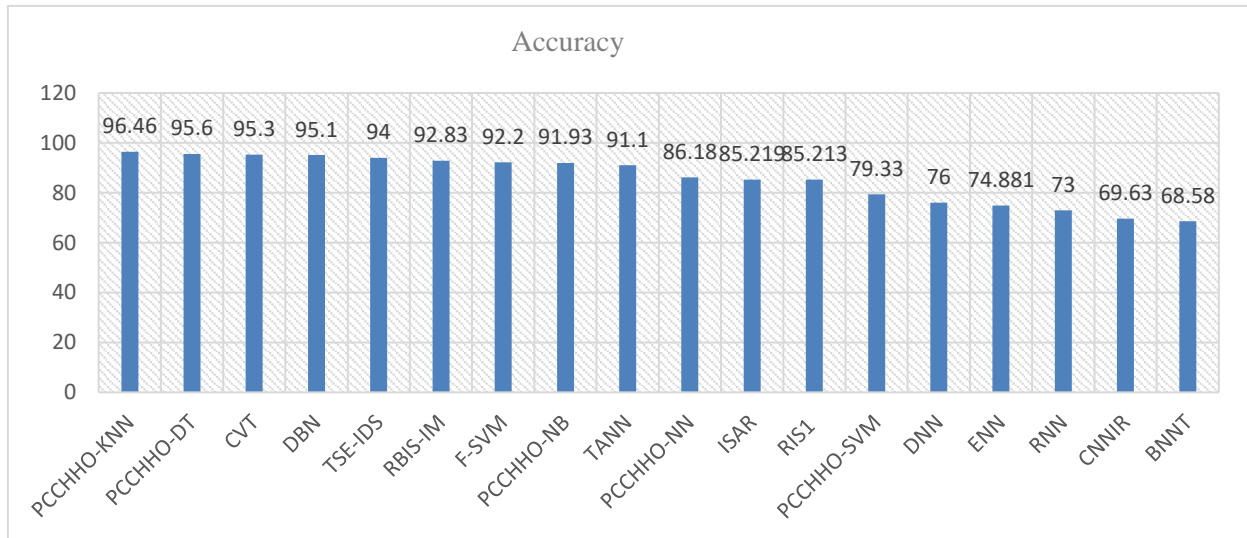


Fig 5. Comparing the accuracy of the proposed method with other methods.



Paper Type (Research paper)

Improving the Performance of Permanent Magnet Synchronous Motor by Using Direct Current Control Method Based On Predictive Controller with Continuous Control Set

Hamid Rahimi Esfahani¹, Reza Sharifian Dastjerdi^{2*}

^{1,2}Department of Electrical Engineering, Lenjan Branch, Islamic Azad University, Isfahan, Iran

Article Info

Article History:

Received: 2024/10/15

Revised: 2024/10/31

Accepted: 2024/11/15

DOI:

Keywords:

Direct current control,
predictive controller,
permanent magnet synchronous
motor drives, total harmonic
distortion fault

*Corresponding Author's Email
Address: Rezasharifian@iau.ac.ir

Abstract

In this paper, direct current control method is proposed to improve the performance of permanent magnet synchronous motor in steady and transient states. In this method, direct flow control based on predictive controller with continuous control set is provided. Thus, In steady state, an active voltage vector along with a zero-voltage vector is applied to the motor in each control cycle. The values of the phase, amplitude and duty factor of the active voltage vector are optimized in such a way that the stator current error is minimized. In the transient mode, to improve the dynamic response of the torque, a voltage vector with maximum amplitude is applied to the motor in the entire control cycle, and the angle of the voltage vector is calculated in such a way that the stator current error is reduced to zero at the end of the control cycle. Spatial vector modulation is used to generate the selected voltage vector. The performance of the method has been evaluated in MATLAB software; The obtained results show that the proposed method of harmonic distortion of the total stator current in the steady state reduces and improves the dynamic response of the motor in the transient state. In addition, the performance of the two presented methods has been compared with a number of recent control methods, and the results show that by using the proposed methods, the performance of the steady and transient state is improved.

1. Introduction

In the past, direct current motors, especially independent excitation direct current motors, were widely used in variable speed drives. But direct current machines had problems such as high price, high weight and volume, presence of commutator, presence of excitation losses, repair and maintenance. With the emergence of new permanent magnet materials with high energy density at the end of the 20th century, a great change was made in the structure of direct current machines. The use of these materials led to the elimination of the coil and the external energy source in the excitation of direct current motors.

On the other hand, progress in the field of high-power semiconductor devices led to the expansion of the use of inverters, and inverters could replace mechanical commutators, and this was the beginning of the construction of permanent magnet synchronous motors and brushless direct current motors. For more than two decades now, permanent magnet synchronous machines have been used in variable speed drives and for normal applications in the power range of several watts to several kilowatts [1].

For the drive of AC motors, the vector control method is actually an evolution of the control

methods created in 1970 by Blash and Hayes as a new control method. If in this method, the size and phase angle of the stator phase currents are controlled in order to create precise control over the motor. Vector control can be implemented in the reference frame corresponding to the coordinates of the stator, rotor or the space vector of the magnetizing flux and using the currents of the stator axes d and q defined in the corresponding coordinates. Generally, in drive systems, vector control is used to control the stator current from the current sensor in order to obtain the actual value of the motor current. Then the measured values are compared with their corresponding reference values. These DC currents are transferred to the reference frame of the three-axis stator to be used by the hysteresis method to generate switching pulses, or to be used by the current controller to generate the required reference voltage in controllers based on SVM¹ [2]. The current vector control method using hysteresis controllers for current control has disadvantages because the inverter's switching frequency is variable and its effective value depends on the motor rotation speed and load torque. Also, the current fluctuations are limited to the bandwidth of the hysteresis controllers. If this bandwidth is large, the harmonics of the stator current increase, and if the bandwidth is small, the switching frequency increases significantly [3-4]. Also, the vector control method based on proportional-integral controllers and SVM modulation in order to control the flow, have a slow dynamic response and high switching frequency, and in this category of controllers, generally, the process of transferring quantities from One reference frame is needed for another reference frame, which significantly increases the dependence of the accuracy of the control method on the motor parameters, and the amount of calculations necessary to implement the control method increases. In order to solve the above problems, predictive controllers can be used [5-6], which are widely used in industrial and commercial applications today. In this category of controllers, a precise mathematical model of the motor is obtained and based on this model, the behavior of the motor is predicted for one to several future cycles. In this category of controllers, an objective function is used to achieve different control goals. The working procedure in these controllers is such that in each

sampling cycle, the behavior of the motor is predicted and based on the predicted value and the reference value, the objective function is evaluated and the voltage vector which has the lowest cost value is calculated as the optimum voltage vector.

There are many types of predictive controllers that differ in the set of available voltage vectors for motor control, the way to choose the voltage vector applied to the motor, and the goals of the control system. For example, in [7], the relationship between torque changes and the voltage vector applied to the motor is obtained, and then the amplitude of the voltage vector and the duration of its application to the motor are assumed to be constant, and the angle of the voltage vector is calculated in such a way that the motor torque at the end of the cycle Sampling is equal to reference torque. In another category of predictive control methods, the state space model of the motor is obtained, and then, using existing control techniques, the optimal voltage vector is selected to meet a specific control goal [8-10]. For example, in [8], a cost function is defined based on the motor state space model. Then the effect of each of the active voltage vectors in this cost function is investigated and the voltage vector that minimizes the cost function is selected as the optimal active voltage vector and is applied to the motor. In [9], after obtaining the system state space model, an active voltage vector and a zero-voltage vector are applied to the motor with the aim of reducing torque fluctuations. Due to the fact that performing the operations related to the selection of the active voltage vector and the calculation of the duty factor in a row leads to the passage of a lot of time, the efficiency of the method decreases at low speeds; Therefore, in this reference, in order to solve this problem, the mentioned operations are performed simultaneously. There is another category of predictive- control methods in which control operations are performed in such a way that the motor flux and torque reach their reference values exactly after one sampling cycle (at the end of each cycle). This control method is called deadbeat control method; This means that the error between the motor flux and torque and their reference values becomes zero exactly after one sampling cycle [11-13]. In predictive controllers, first, the mathematical model of the machine is obtained in the static reference frame, and based on that, the components of the stator current are estimated for the next sampling cycle. Then the RMS function of the stator current error is

¹ Space Vector Modulation

obtained and the parameters of the voltage vector applied to the motor are optimized in such a way that the minimum effective value of the error for the stator current components is obtained. By using these controllers, the stator current harmonics are significantly reduced. In the MPDCC² methods presented in recent years, modulation blocks are usually not used and only a limited number of voltage vectors have been available to apply to the motor. In some of these methods, only the duration of applying the selected voltage vector to the motor is optimized, which is much less effective compared to the method presented in this research, because in the proposed method in this research, not only the duration of applying the vector voltage to the motor is optimized, but two other parameters of the voltage vector, i.e. its amplitude and phase, are also optimized; Therefore, by using the proposed method in this project, the stator current harmonics are significantly reduced. The structure of the paper according to the explanations given in this part is such that the proposed direct flow control methods to improve the direct flow control have been evaluated in the second part and transient and dynamic stability of the motor are included in the proposed method. In the third part, the results of the simulations are shown in such a way that the efficiency of the proposed method is clearly shown. At the end, conclusions and suggestions for further work are given in the fourth part.

2. The Direct Control Method of the Proposed Predicted Current

The best performance of a control system is achieved when all the parameters of the voltage vector, i.e. the range, phase and duration of applying the voltage vector to the motor are optimized by the control system. Therefore, in order to effectively reduce stator current harmonics, torque fluctuations and stator flux fluctuations, in this paper, a continuous control set - predictive model current controller (CCS-MPDCC³) is presented, by which all parameters of the voltage vector are optimized. In this part, first, the features of predictive controllers based on system modeling for flow control will be explained, and then the direct flow control method presented in this paper will be explained in detail.

The very important feature of the MPC⁴ controller, which has turned it into a very powerful tool in control, is that nonlinear systems with multiple outputs and multiple inputs can be easily modeled and controlled by this control system.

A) Control System Model

Usually, MPC controllers are modeled in discrete space with fixed sampling time T_s . In this way, the inputs of the system are necessarily changed in the moments that are integer coefficients of T_s ; In other words, at moments $t = kT_s$ where $k \in \{0,1,2, \dots\}$ is and it indicates the number of samples. Because power electronic applications generally have nonlinear dynamics, it is more common to model the controlled system in the form of equation 1 in the nonlinear state space, where $x(k)$ represents the value of the state variable in the k th sample and $u(k)$ represents the input value of the system in the k th sample:

$$x(k+1) = f(x(k), u(k)), \quad k \in \{0,1,2, \dots\} \quad (1)$$

As mentioned, in this part, MPC controllers with continuous output are examined, in which the output of the control system is given to the modulator unit, and then the status of the power switches is determined by the modulator. Therefore, the constraint related to system inputs can be defined according to equation (2).

$$u(k) \in U \subseteq \mathbb{R}^p, \quad k \in \{0,1,2, \dots\} \quad (2)$$

where p is the number of system keys and the set U is determined according to the type of inputs. The system input $u(k)$ can be a voltage vector, or a duty factor, or the state of a key. For example, if $u(k)$ is the state of the keys, the set U is equal to $U = [0, 1]^p$ [14]. In addition to the constraints in the inputs can be considered by MPC, the constraints in the controlled states can also be considered. For example:

$$x(k) \in X \subseteq \mathbb{R}^n, \quad k \in \{0,1,2, \dots\} \quad (3)$$

where n is the number of state variables. For example, the constraint related to the state can be the voltage of a capacitor in a converter, or the voltage of the neutral point in a three-level inverter, or the inductor current in a resistive-inductive load.

b) Cost Function

In the MPC method, in each sampling cycle for a specific state $x(k)$ (obtained by estimation or measurement) and for several dimension N samples, the pre-defined cost function related to

² Model predictive Direct Current Control

³ Continious Control Set-Model predictive Direct Current Control

⁴ Model predictive Control

the control objectives is minimized. The cost function introduced in (4) is a general form of the cost function that has been widely used in recent years' papers.

$$V(x(k), \vec{u}'(k)) \triangleq F(x'(k+N)) + \sum_{l=k}^{k+N-1} L(x'(l), u'(l)) \quad (4)$$

In relation (4), $L(\dots)$ and $F(\dots)$ are weighting functions that are used to determine the appropriateness of the system's behavior for a specific input; For example, they specify the amount of error between the reference voltage value and the predicted value. The predicted values related to the system state are formed according to equation (5):

$$x'(l+1) = f(x'(l), u'(l)), \quad l \in \{k, k+1, k+2, \dots, k+N-1\} \quad (5)$$

in which

$$u'(l) \in \mathbb{U}, l \in \{k, k+1, k+2, \dots, k+N-1\} \quad (6)$$

They represent the experimental inputs of the control system. The recursive relationship stated in (5) is initialized with the current value of the system states. It means:

$$x'(k) \leftarrow x(k) \quad (7)$$

Therefore, equation (5) is the predictive expression of the state of the control system, which is obtained by applying the inputs stated in (8) in $\{k, k+1, k+2, \dots, k+N-1\}$ samples.

$$\vec{u}'(k) = \{u'(k), u'(k+1), \dots, u'(k+N-1)\} \quad (8)$$

Predicted state variables and system inputs are both bound according to equation (2), so we have:

$$u'(l) \in \mathbb{U}, \quad x'(l) \in \mathbb{X} \quad \forall l \in \{k, k+1, k+2, \dots, k+N-1\} \quad (9)$$

In addition, it is usually necessary that $x'(k+N)$ satisfy a certain condition constraint; For example $x'(k+N) \in \mathbb{X}$

It is necessary to pay attention to this point that the selection of the constraint related to the state $x'(k+N)$ is usually done according to the stability problem. According to the above explanations, the constrained optimization introduced in relation (4) determines the order of the system inputs to realize the optimal control in the k th sample and for the state $x(k)$.

$$\vec{u}'(k) \triangleq \{u(k;k), u(k+1;k), \dots, u(k+N-1;k)\} \quad (11)$$

c) How to Solve the Optimization Problem and Select the Voltage Vector with the Passing of Sampling Cycles

Despite the fact that $\vec{u}'(k)$ has the optimal voltage vector for all future cycles, only the first voltage vector is selected by the control system and applied to the motor. That is, the input of the

control system is adjusted according to equation (12):

$$u(k) \leftarrow u(k;k) \quad (12)$$

In the next sampling cycle, i.e. the $k+1$ th sample, the system states $x(k+1)$ are measured (or estimated), then all the above actions are repeated again for the new cycle and the optimal voltage vector $\vec{u}'(k+1)$ is obtained. If, this set includes the optimal voltage for all future cycles, but only its initial voltage vector is selected to be applied to the motor, i.e. $u(k+1) \leftarrow u(k+1;k+1)$. This process is repeated for all subsequent cycles. To clarify the meaning of the mentioned content, Figure(1) shows how to perform operations during the sampling cycles and select the optimal voltage vector for the case where the predictive controller has N equal to 3. Therefore, the MPC method can be called an open-loop optimal method.

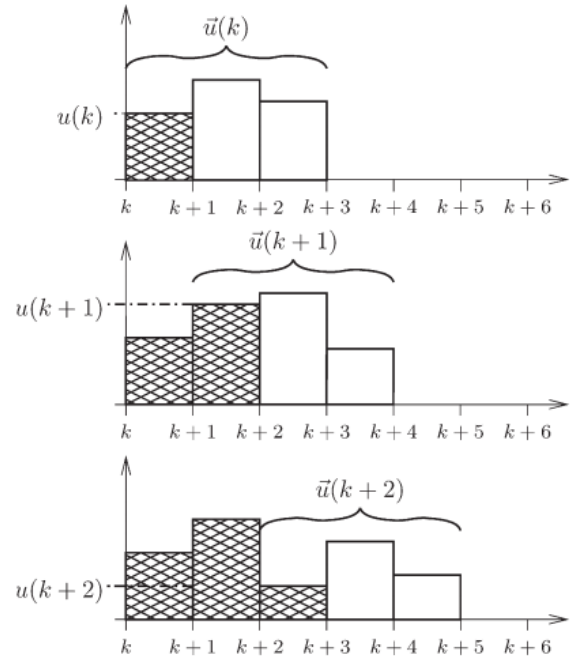


Figure (1): How to perform operations for the case where N is equal to 3.

d) Parameters Related to Design

As seen in the previous section, the MPC controller has the ability to control a nonlinear multivariable system. In order to implement the predictive controller, in addition to choosing the appropriate value for the duration of a sampling cycle, i.e. T_s , it is necessary to choose the appropriate objective function to achieve the control goals. In other words, suitable functions for weighting functions, ie, $L(\dots)$ and $F(\dots)$ should be selected. Also, the number of cycles that need to be predicted in each cycle (horizon), i.e. N ,

should be selected based on control objectives. If possible, the constraints related to the inputs and states of the system, i.e. X and X_f , can also be designed.

3. The Controller of the Continuous Control Set, the Direct Control of the Flow of the Proposed Prediction Model

In the proposed CCS-MPDCC method, the stator current components are the controlled variables and spatial vector modulation has been used to produce the selected optimal voltage vector. Therefore, due to the use of SVM, the amplitude and phase of the voltage vector can be adjusted to any desired value based on the objectives of the controller. The main purpose of this method is to minimize stator current ripples in steady state and obtain a fast- dynamic response for motor torque in transient state. The process of obtaining the optimal voltage vector parameters is shown in detail in the following sections:

3-1) Improved Steady State Performance

a) Permanent Magnet Synchronous Motor Model

In this proposed method, in order to prevent the transfer of rotating coordinates and reduce the complexity of calculations, the static reference frame is used to formulate the motor equations, the motor equations can be considered from [13].

b) Changes in Stator Current Components

Since the proposed method is implemented using a digital processor, it is necessary to obtain the discrete time representation of equations (16) and (17). If by discretizing them, we have:

$$\Delta i_{sd} = \frac{T_s}{L_d} (u_{sd} - Ri_{sd} + \omega_r \psi_f \sin \theta_r) \quad (13)$$

$$\Delta i_{sq} = \frac{T_s}{L_q} (u_{sq} - Ri_{sq} - \omega_r \psi_f \cos \theta_r) \quad (14)$$

According to the above equations, the values of i_{sd} , i_{sq} , ω_r and θ_r can be calculated in each control cycle, by proper control of the voltage vector applied to the motor, i.e. u_{sd} and u_{sq} in equations (18) and (19), controlling the stator current components towards Their reference values become possible. If the following parameters are defined:

- The length of the control cycle is T_s .
- An arbitrary non-zero voltage vector $(\vec{u}_s^* = [u_{sd}^* \ u_{sq}^*]^T)$ with duration T_k^* ($T_k^* \leq T_s$) is applied to the motor.

- A zero voltage vector (ZVV) that is applied to the motor in the remainder of the control cycle.

Considering these definitions and based on the variables of equations (13) and (14), the stator current components due to the application of a ZVV and the reference voltage vector (\vec{u}_s^*) are formulated as equations (15) to (18):

$$\Delta i_{sd0} = \frac{(T_s - T_k^*)}{L_d} (-Ri_{sd} + \omega_r \psi_f \sin \theta_r) \triangleq S_{d0}(T_s - T_k^*) \quad (15)$$

$$\Delta i_{sq0} = \frac{(T_s - T_k^*)}{L_q} (-Ri_{sq} - \omega_r \psi_f \cos \theta_r) \triangleq S_{q0}(T_s - T_k^*) \quad (16)$$

$$\Delta i_{sd1} = \frac{T_k^*}{L_d} (u_{sd}^* - Ri_{sd} + \omega_r \psi_f \sin \theta_r) \triangleq \frac{T_k^*}{L_d} u_{sd}^* + S_{d0} T_k^* \quad (17)$$

$$\Delta i_{sq1} = \frac{T_k^*}{L_q} (u_{sq}^* - Ri_{sq} - \omega_r \psi_f \cos \theta_r) \triangleq \frac{T_k^*}{L_q} u_{sq}^* + S_{q0} T_k^* \quad (18)$$

if

$$S_{d0} = \frac{1}{L_d} (-Ri_{sd} + \omega_r \psi_f \sin \theta_r) \quad (19)$$

$$S_{q0} = \frac{1}{L_q} (-Ri_{sq} - \omega_r \psi_f \cos \theta_r) \quad (20)$$

In equations (15) to (18), Δi_{sd0} and Δi_{sd1} are the changes of stator current components d, q respectively, due to the application of a ZVV and \vec{u}_s^* for the motor. If T_s is included in about 100 microseconds [15-16], assuming that all variables are constant in a control cycle; Therefore, the slopes of S_{d0} and S_{q0} are the same in a control cycle. As a result, the value of the stator current components at the end of the Kth control cycle (or the beginning of the control cycle (k+1) are calculated as follows:

$$i_{sd}(k+1) = i_{sd}(k) + \frac{T_k^*}{L_d} u_{sd}^* + S_{d0} T_s \quad (21)$$

$$i_{sq}(k+1) = i_{sq}(k) + \frac{T_k^*}{L_q} u_{sq}^* + S_{q0} T_s \quad (22)$$

which, $i_{sd}(k)$ are the d axis stator current at the beginning of the kth control cycle.

c) Minimizing Stator Current Ripples

In general, in order to evaluate the performance of a signal that is different from the reference signal, the square root is a common measurement [17]. The RMS value of the stator current error on a control cycle is defined as follows:

$$|\vec{i}_{s-error(RMS)}|^2 = \frac{1}{T_s} \int_0^{T_s} \{ (i_{sd}^* - i_{sd}(k+1))^2 + (i_{sq}^* - i_{sq}(k+1))^2 \} dt \quad (23)$$

i_{sd}^* and i_{sq}^* are the reference values of stator current components. Using equations (26) and (27) and simplifications of equation (28), the following equation is obtained.

$$\left| \vec{i}_{s-error(RMS)} \right|^2 = \frac{1}{T_s} \int_0^{T_s} \left\{ (i_{sd-error} - \frac{T_k^*}{L_d} u_{sd}^* - S_{d0} T_s)^2 + (i_{sq-error} - \frac{T_k^*}{L_q} u_{sq}^* - S_{q0} T_s)^2 \right\} dt \quad (24)$$

If

$$i_{sd-error} = i_{sd}^* - i_{sd}(k) \quad (25)$$

$$i_{sq-error} = i_{sq}^* - i_{sq}(k) \quad (26)$$

The function $\left| \vec{i}_{s-error(RMS)} \right|^2$ in equation (24) is a cost function, whose variables are T_k^* , u_{sd}^* and u_{sq}^* . In order to obtain the minimum of stator current ripples, the cost function optimization problem must be solved. The form of this problem is similar to the quadratic optimization problem (QP) in references [18]. As a result, this optimization problem is highlighted and will have an optimal solution in the practical area [19]. By minimizing the optimization problem, the equations of the components of the optimal voltage vector and the optimal duration are given as the following equations.

$$u_{sd}^* = \frac{L_d(i_{sd-error} - S_{d0} T_s)}{k T_s} \quad (27)$$

$$u_{sq}^* = \frac{L_q(i_{sq-error} - S_{q0} T_s)}{k T_s} \quad (28)$$

$$T_k^* = \frac{i_{sq-error} - S_{q0} T_s}{2S_{q1} - S_{q0}} + \frac{i_{sd-error} - S_{d0} T_s}{2S_{d1} - S_{d0}} \quad (29)$$

After calculating u_{sd}^* , u_{sq}^* and T_k^* , these values are applied to the SVM block. The SVM block generates two adjacent active voltage vectors (AVV) based on the values of u_{sd}^* and u_{sq}^* , and then they are applied to the motor in the time period calculated according to the values of u_{sd}^* , u_{sq}^* and T_k^* , are used. In summary, based on the above, if the d-axis and q-axis components, the applied voltage vector of the motor, are adjusted according to equations (27) and (28), and this applied voltage vector for the motor is calculated in a period of time according to equation (29), the ripples The minimum stator current is obtained in steady state.

d) Reducing the Switching Frequency of the Inverter

According to the explanations given, in steady state, two adjacent AVVs (which are selected based on the components of the optimal voltage vector) with a zero-voltage vector (ZVV) to the motor respectively for time intervals T_k^* and $(T_s - T_k^*)$ is applied. According to reference [17], ZVV should be applied symmetrically at the beginning and end of the control cycle to obtain

the minimum RMS ripple. Since two-level VSI⁵ is used to implement the proposed method, ZVV can be u_0 or u_7 . To reduce the switching frequency, ZVV should be selected based on the adjacent voltage vectors in such a way that the minimum switching transfer is obtained. It means that if SVM uses u_0 or u_7 as ZVV to combine with the selected voltage vector, then u_0 or u_7 must be selected by the control system. By applying this method, only the state of one switch is changed at each moment of switching. Therefore, unnecessary switching transfers are prevented and switching losses are reduced.

e) Improving Motor Efficiency Using the Principle of Maximum Torque in Terms of Amperes

One of the ways to improve motor efficiency is to use the principle of maximum torque in terms of amperes in the range below the rated speed. In this method, by optimally adjusting the magnetic flux of the stator, the losses are significantly reduced and, as a result, the efficiency is increased. In this research, using this principle, the current components are calculated in different working conditions and stored in the memory of the digital processor in the form of an observation table. During the implementation of the control algorithm, according to the working conditions of the motor, the stator current components are fetched from this table.

3-2) Improving Transient Mode Performance

Improving the dynamic response of the motor is the main goal of the controller in transient mode. Therefore, the parameters of the voltage vector (phase, amplitude and duration) should be set in a different state of the proposed method from the previous parts of the steady state. In the transient state, the parameters of the voltage vector are adjusted so that the actual value of the stator current components reach their reference values at the end of the control cycle. It means that the deadbeat control method, which has a fast-dynamic response, should be applied to control the stator current components in the transient state. As a result, the following set of equations should be considered to find the optimal voltage vector parameters:

⁵ Voltage Source Inverter

$$\Rightarrow \begin{cases} i_{sd}(k+1) = i_{sd}^* \\ i_{sq}(k+1) = i_{sq}^* \\ i_{sd}(k) + \frac{T_k^*}{L_d} u_{sd}^* + S_{d0} T_s = i_{sd}^* \\ i_{sq}(k) + \frac{T_k^*}{L_q} u_{sq}^* + S_{q0} T_s = i_{sq}^* \end{cases} \quad (30)$$

In order to obtain a fast-dynamic response in the transient state, a voltage vector with the largest amplitude is applied to the motor in the entire control cycle. It means that, in the transient state, T_k^* and $|\vec{u}_s^*|$ are set equal to T_s and V_{max} , respectively, where V_{max} is the maximum voltage range that can be obtained in the linear region of SVM. Therefore, by replacing T_k^* and $|\vec{u}_s^*|$ in equation (30) and simplifying, the optimal voltage vector phase equation is obtained based on equation (31):

$$\alpha_s^* = \tan^{-1} \left(\frac{L_q(i_{sq-error} - S_{q0}T_s)}{L_d(i_{sd-error} - S_{d0}T_s)} \right) \quad (31)$$

where α_s is the phase of the voltage vector applied to the motor relative to the axis d of the stationary reference frame. In summary, based on the above information, if the amplitude of the stator voltage vector V_{max} is adjusted and the phase of the voltage vector is calculated based on equation 36 and this voltage vector is applied in all motor control cycles, it effectively improves the dynamic response of the motor in the transient state. It should be noted that in the implementation of the deadbeat control method, the current error is measured at the beginning of each control cycle and must be zero at the end of each control cycle, so no excess current occurs in the transient state.

3-3) Identifying Steady State and Transient State

After calculating the components of the optimal voltage vector from equations (32) and (33), the amplitude of the voltage vector is calculated from equation (32). If the harmonics of the stator current in the over modulation region of SVM are higher than the harmonics in the linear region. Over modulation area is not used in the proposed method. Therefore, the diagnosis of steady state and transient state is done on the basis that if $\frac{T_k^*}{T_s}$ is greater than one and $|\vec{u}_s^*|$ is greater than V_{max} , it is not possible to combine the optimal voltage vector in a control cycle used in the linear region of SVM. In these conditions, the range of the voltage vector is limited to V_{max} and T_k^* is limited to T_s . Therefore, the stable state can be identified

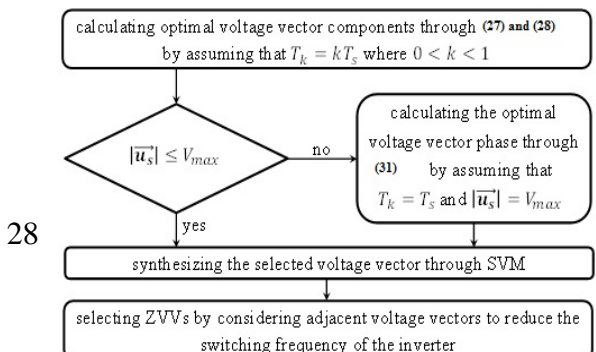
by having the conditions of equations (33) and (34) and otherwise it is a transient state:

$$|\vec{u}_s^*| = \sqrt{(u_{sd}^*)^2 + (u_{sq}^*)^2} \quad (32)$$

$$\frac{T_k^*}{T_s} \leq 1 \quad (33)$$

$$|\vec{u}_s^*| < V_{max} \quad (34)$$

Based on the explanations provided in this section, the calculation steps of the proposed CCS-MPCC method in each control cycle are summarized in Figure (2). In this research, in order to effectively reduce stator current harmonics, which leads to reduction of harmonic losses, reduction of torque fluctuations and reduction of stator flux fluctuations, direct current control method based on predictive controller with continuous control set (CCS- MPDCC) is used. In this method, in order to achieve the optimal voltage vector with the desired amplitude and phase, SVM spatial vector modulation is used, and all parameters of the voltage vector are optimized with the aim of minimizing stator current harmonics. In the presented controller, first, the mathematical model of the machine is obtained in the static reference frame, and based on that, the stator current components are estimated for the next sampling cycle. Then the RMS function of the stator current error is obtained and the parameters of the voltage vector applied to the motor are optimized in such a way that the minimum RMS value of the error for the stator current components is obtained. By using this controller, the harmonics of the stator current are significantly reduced because, firstly, the RMS function of the error, which is the best indicator for evaluating the tracking quality of the reference signal by another signal, has been selected as the target function. And secondly, all parameters of the voltage vector are optimized simultaneously. Thirdly, the SVM modulation method is used in order to reach the selected optimal voltage vector. In the MPDCC methods presented in recent years, modulation blocks were not used and only a limited number of voltage vectors were available to be applied to the motor. In some of these methods, only the duration of applying the selected voltage vector to the motor is optimized, which is much less effective compared to the method presented in this research, because in the proposed method in this project, not only the duration of applying the voltage vector to the



motor is optimized, but two other parameters of the voltage vector, i.e. its amplitude and phase, are also optimized; Therefore, by using the proposed method in this research, the stator current harmonics are significantly reduced.

Figure (2): Calculation steps of the proposed CCS-MPCC method

4) Results of Simulations Related to the Presented Method

In order to evaluate the efficiency of the methods presented in this research, simulation was used in MATLAB software. If during simulation in the MATLAB/Simulink software environment, it is possible to view all motor quantities online. In such a way that all the information to be studied in each sampling cycle (for example, once every 100 microseconds) is calculated by the control system and then various processing operations are performed on the information by MATLAB software; For example, the harmonic spectrum of the stator phase currents is obtained and the graph of instantaneous changes of all quantities is drawn and the effective switching frequency is also calculated. The effectiveness of the proposed CCS-MPCC method is shown in this section using MATLAB/Simulink software. The implementation of the proposed method has been done carefully and then it has been compared with the conventional hysteresis implementation based on the MPCC method and the method proposed in reference [20]. In order to obtain minimum stator current ripples in the hysteresis method based on the MPCC method, the bandwidth of the two hysteresis controllers is set to 0.07 [21]. Motor and control system parameters in all simulations and practical results are shown in Table (1). As can be seen, the sampling time of all methods is set to 100 microseconds according to the sampling time of the methods implemented in references [22-26]. To simplify the continuation of this research, the hysteresis-based DCC⁶ method, the proposed MPCC method in reference [20], and the proposed CCS-MPCC method are given with DCC, Duty-MPCC, and CCS-MPCC, respectively.

Table (1): control system and motor parameters

Parameters	Value
Number of pole pairs	2
Permanent-magnet flux	0.901 (Wb)
Stator resistance	7.9 (Ω)
d -axis stator inductance	0.070 (H)
q -axis stator inductance	0.117 (H)
Rated speed	1500 (rpm)
Rated voltage	400 (V)
Rated torque	1.4 (Nm)
Rated current	0.7 (A)
Rated power	200 (W)
DC link voltage	250 (V)
Sampling time (T_s)	100 (μs)
Bandwidth of the hysteresis current controller in DCC method	0.07

a) Harmonic spectrum of stator current

In order to show the performance of the proposed steady state method, the harmonic spectrum of one phase of the stator current at 100% of the nominal speed is drawn in Figure (3). that the rated load is applied to the motor. In this simulation, the THD⁷ of the stator current has been calculated up to 6 kHz. Based on these figures, it is clear that in all speed ranges, the lowest THD of the stator current is obtained by applying the proposed CCS-MPCC method, because in this method, the parameters of the voltage vector are optimized with the aim of minimizing the stator current ripples. have been made Based on Figure 3 to 5, the numerical comparison of THDs at rated speed shows that a

⁶ Direct Current Control

⁷ Total Harmonic Distortion

73.26 percent THD reduction is achieved by applying the proposed CCS-MPCC method compared to the DCC method. While the reduction of THD, by applying Duty-MPCC methods, is 66.12%. The simulated results in the figure confirm that THD is effectively reduced by applying the proposed method.

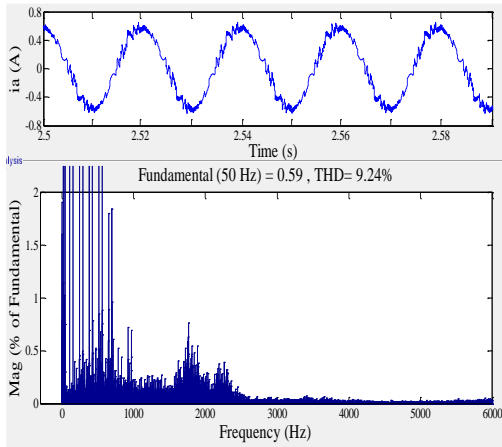


Figure (3): Harmonic spectrum at 1500 rpm DCC method

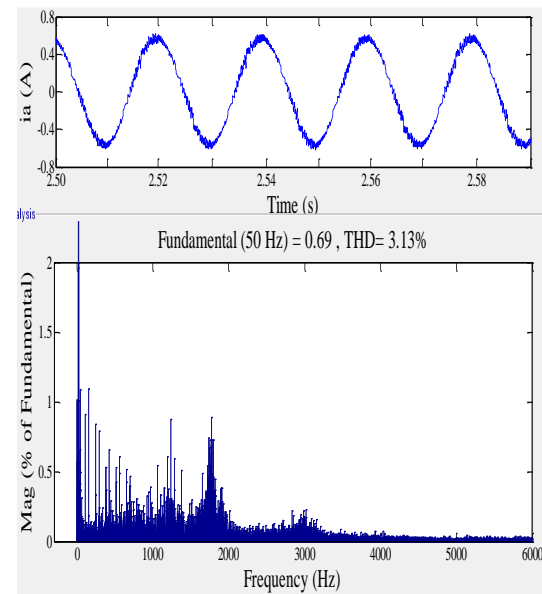


Figure (4): Harmonic spectrum at 1500 rpm Duty-MPCC method

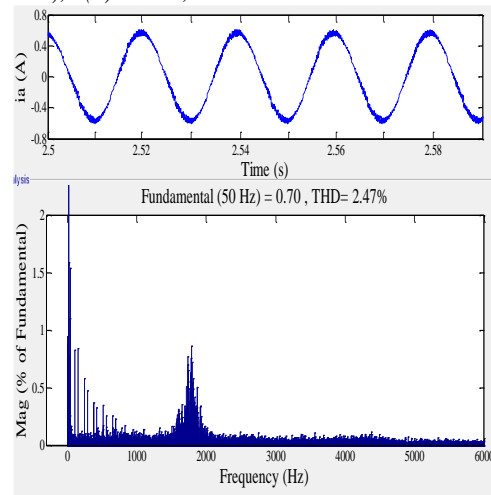


Figure (5): Harmonic spectrum at 1500 rpm Proposed CCS-MPCC method

b) Effect of Average Switching Frequency

The average switching frequency is one of the important issues in comparing the performance of the stable modes of the methods. It is clearly seen that a higher (lower) average switching frequency leads to a lower (higher) current ripple and as a result, it will have a lower (higher) THD. Therefore, in order to have a specific comparison, the average switching frequency of all methods should be considered the same. It should be noted that the average switching frequency of the inverter is calculated based on equation 35, if N is the number of switching transitions of the inverter during a constant period, for example, in the period of 200 milliseconds in this experiment; and K is the number of inverter switches, which is equal to 6 in two-level VSI [27-28].

$$f_{ave} = N/K/\Delta t \quad (35)$$

To investigate the effect of average switching frequency on the implementation of the DCC method, a similar simulation as before has been carried out under the same conditions but with a sampling time of 40 microseconds (that is, a sampling frequency of 25 kHz). In DCC, when the sampling frequency increases to 25 kHz, the switching frequency increases to about 7.5 kHz. By increasing the switching frequency, THD reduction results, and choosing an optimal voltage vector with optimal parameters can further reduce THD. Apart from this, in order to implement the DCC method at a sampling frequency of 25 kHz, more powerful digital processors and faster analog-to-digital converters are needed, which significantly increases the implementation costs. In Table (2), the THD value in all methods is summarized.

Table (2): Comparison of THD values at a speed of 1500 rpm

Method	THD	f_{ave} (kHz)
DCC (at 10-kHz)	9.24 %	3.1 %
duty-MPCC	3.13 %	7.5 %
proposed CCS-MPCC	2.47 %	7.5 %
DCC (at 25-kHz)	3.70 %	7.5 %

c) Startup response

In order to test the dynamic response of the torque in all methods, the motor starting response is checked from the stop state to the rated speed in no-load conditions. As in all methods, PI controllers with specific integral and proportional gains have been used to regulate the motor speed. The gains of the PI controller are adjusted so that the motor speed is controlled without any overshoot. In such a way that the starting response of these methods is almost the same and there is no steady state error in the speed and torque of the motor. In all methods, after starting the motor, the speed reaches its nominal value after 0.33 seconds without any overshoot. For greater clarity of the dynamic torque response comparison, the torque response is enlarged in Figure 6 to 8. If, the torque responses are almost the same in all methods, but there is a very small amount of difference in the time of torque increase.

According to Figures (6) to(8), the dynamic response of the torque in the proposed CCS-MPCC method is faster than the other two methods, because a voltage vector with the optimal amplitude and phase is applied to the motor in the entire control cycle. While in the DCC method, there is no prediction and optimization to obtain a fast torque response, and the torque increase time in the DCC method is 7.9 milliseconds, in the Duty-MPCC method, 8.6 milliseconds and in the proposed CCS-MPCC method is 6.2 milliseconds.

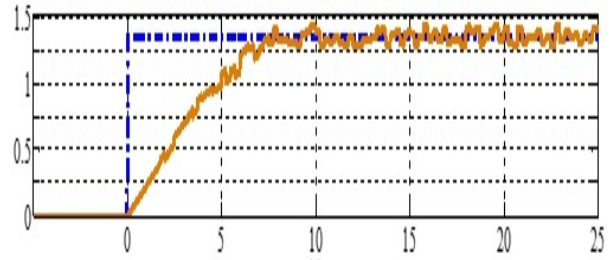


Figure (6): Magnified torque response in start-up test DCC method

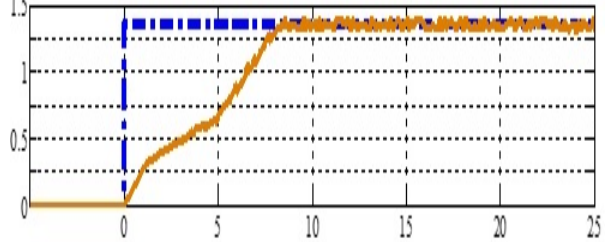


Figure (7): Magnified torque response in start-up test Duty-MPCC method

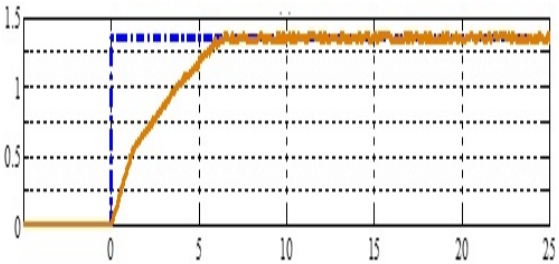


Figure (8): Magnified torque response in start-up test Proposed CCS-MPCC method

It results that the dynamic response of the DCC method is improved by applying the proposed CCS-MPCC method. If the proposed CCS-MPCC method in this research minimizes the stator current error by using SVM space vector modulation. Also, the voltage vector parameters are calculated in each control cycle to find the optimal voltage vector.

5) “Comparative Analysis of the Proposed CCS-MPCC Method with Recent Similar Methods”

This section demonstrates the advantages of the proposed CCS-MPCC method over several recent methods. These advantages are theoretical and will be validated in subsequent sections through simulations and experimental results.

a) Comparison with the Method Presented in Reference [28]

In reference [28], an optimal voltage vector is selected by the control system for simultaneous

control of torque and stator flux. This voltage vector is then applied to the motor using space vector modulation. In this method, the RMS torque ripple is considered as the cost function, while deadbeat control of the stator flux magnitude is considered as a constraint. The Lagrange multiplier method is used to solve the optimization problem and obtain the voltage vector parameters. By comparing the analysis of the proposed method in this reference with the proposed CCS-MPCC method, the following results are obtained:

1) The computational complexity of the proposed CCS-MPCC method is lower than the method presented in this reference.

This is because in reference [28], a closed-loop estimator or observer is required to estimate torque and stator flux in each control cycle, whereas in the proposed CCS-MPCC method, no estimation is required.

2) The voltage vector selection process in this reference is more complex than the process used in the proposed CCS-MPCC method. This is because in reference [28], a constrained optimization problem must be solved to obtain the optimal voltage vector, while in the proposed CCS-MPCC method, an unconstrained optimization problem must be solved.

3) The stator current THD in the proposed CCS-MPCC method is lower than the method presented in this reference.

This is because in the proposed CCS-MPCC method, the stator current components are directly controlled, while in reference [28], torque and stator flux are directly controlled. In other words, the stator current is controlled indirectly. In other words, the stator flux ripple and torque ripple in this reference are less than the proposed CCS-MPCC method.

b) Comparison with the Model Predictive Current Control Method Presented in Reference [20]

In the method proposed in reference [20], an active voltage vector and a zero voltage vector are applied to the motor in each control cycle. In this method, two ways are shown to select the best active voltage vector in each control cycle. In the first way, the cost of each available voltage vector is calculated, and the voltage vector with the lowest cost is selected as the best voltage vector. After selecting the best voltage vector, the deadbeat control method is used to calculate the duration of the selected active voltage vector. In the second way, to avoid checking all available vectors

and selecting the optimal vector, a theoretical optimal voltage vector is calculated for using the deadbeat control method, and then the closest active voltage vector to the optimal voltage vector is selected as the best voltage vector. Finally, the duration of the selected voltage vector is calculated with the aim of minimizing the error between the selected voltage vector and the theoretical optimal voltage vector. As can be seen, in both methods shown in this reference, only one active voltage vector is used for the motor; therefore, the stator current error cannot be minimized.

To overcome these shortcomings, the proposed CCS-MPCC method in this study minimizes the stator current error using space vector modulation (SVM). In the first method shown in this reference, the cost of all active voltage vectors in each control cycle is calculated for the best voltage vector. While in the proposed method, the voltage vector parameters are calculated in each control cycle to find the optimal voltage vector. Also, in this reference, the methods used to control the motor in steady-state are used without any changes for controlling the motor in transient state as well. Whereas, in the proposed CCS-MPCC method in this study, the predictive control model is adapted to minimize the stator current error in steady-state and the deadbeat control method is used to obtain a fast dynamic response in transient state.

6) Conclusion

In the proposed CCS-MPCC method, the components of the voltage vector and the duration of the application of the motor voltage vector are optimized so that the stator current ripples are minimized in the steady state and in the transient state, the largest voltage vector is applied to the motor in the entire control cycle. In this case, in order to control the stator current components in the deadbeat state, the phase of the voltage vector is adjusted so that the error of the stator current components is reduced to zero at the end of the control cycle. The performances of both steady state and transient state of the proposed method have been shown using simulations and practical results, as the following results have been obtained:

The torque ripple, stator flux ripple, and stator current THD of the proposed method are significantly lower than conventional DCC methods, and the dynamic response of the motor in the proposed method is faster than DCC methods.

In summary, the proposed method effectively improves the steady state and transient performances. Therefore, the proposed MPCC method can be considered as a useful algorithm in high-performance PMSM drives that require precise motor control.

References

- [1] D. Casadei, F. Profumo, G. Serra, and A. Tani, "FOC and DTC: Two variable schemes for induction motors torque control," *IEEE Trans. Power Electron.*, vol. 17, no. 5, pp. 779–787, Sep. 2002.
- [2] S. A. Zaid, O. A. Mahgoub, and K. A. El-Metwally, "Implementation of a new fast direct torque control algorithm for induction motor drives," *IET Elec. Power Appl.*, vol. 4, no. 5, pp. 305-313, Mar. 2009.
- [3] M. Paicu, I. Boldea, G. Andreescu, and F. Blaabjerg, "Very low speed performance of active flux based sensorless control: interior permanent magnet synchronous motor vector control versus direct torque and flux control," *IET Power Electron.*, vol. 3, no. 6, pp. 551-561, Nov. 2009.
- [4] B. Singh, S. Jain, and S. Dwivedi, "Torque ripple reduction technique with improved flux response for a direct torque control induction motor drive," *IET Power Electron.*, vol. 6, no. 2, pp. 326-342, Jun. 2013.
- [5] Y. Inoue, S. Morimoto, and M. Sanada, "A novel control scheme for maximum power operation of synchronous reluctance motors including maximum torque per flux control," *IEEE Trans. Ind. Appl.*, vol. 47, no. 1, pp. 115-121, Jan. 2011.
- [6] G. Heins, M. Thiele, and T. Brown, "Accurate torque ripple measurement for PMSM" *IEEE Trans. Instrum. Meas.*, vol. 60, no. 12, pp. 3868-3874, Nov. 2011.
- [7] H. Zhu, X. Xiao and Y. Li, "Torque ripple Reduction of the torque predictive control scheme for permanent-magnet synchronous motors," *IEEE Trans. Ind. Electron.*, vol. 59, no. 2, pp. 871-877, Oct. 2012.
- [8] R. Morales-Caporal, and M. Pacas, "Encoderless Predictive direct torque control for synchronous Reluctance machines at very low and zero speed," *IEEE Trans. Ind. Electron.*, vol. 55, no. 12, pp. 4408-4416, Dec. 2008.
- [9] T. Geyer, and S. Mastellone, "Model predictive direct torque control of a five-level ANPC converter drive system," *IEEE Trans. Ind. Appl.*, vol. 48, no. 5, pp. 1565-1575, Oct. 2012.
- [10] T. Burtscher, and T. Geyer, "Deadlock Avoidance in Model predictive direct torque control," *IEEE Trans. Ind. Appl.*, vol. 49, no. 5, pp. 2126-2135, Sep. 2013.
- [11] W. Song, J. Ma, L. Zhou, and X. Feng, "Deadbeat predictive power control of single-phase three-level neutral-point-clamped converters using space-vector modulation for electrical railway traction" *IEEE Trans. Power Electron.*, vol. 31, no. 1, pp. 721-732, Sep. 2015.
- [12] Y. Wang, T. Ito, and R. D. Lorenz, "Loss manipulation capabilities of deadbeat direct torque and flux control induction machine drives," *IEEE Trans. Ind. Appl.*, Accepted for publication in 2015.
- [13] R. Sh. Dastjerdi, M. A. Abasian, and etc., "Performance Improvement of Permanent-Magnet Synchronous Motor Using a New Deadbeat-Direct Current Controller," *IEEE Trans. Ind. Appl.*, vol. 50, no. 3, pp. 3530-3543, May. 2018.
- [14] S. E. Dreyfus, "Some types of optimal control of stochastic systems," *SIAM J. Control Opt.*, vol. 2, no. 1, pp. 120–134, Jan. 1964.
- [15] G. Abad, M. A. Rodriguez, and J. Poza, "Two-level VSC based predictive direct torque control of the doubly fed induction machine with reduced torque and flux ripples at low constant switching frequency," *IEEE Trans. Power Electron.*, vol. 23, no. 3, pp. 1050–1061, May 2008.
- [16] F. Niu, K. Li, and Y. Wang, "Direct torque control For permanent magnet synchronous machines based on duty ratio modulation," *IEEE Trans. Ind. Electron.*, vol. 62, no. 10, pp. 6160-6170, Oct. 2015.
- [17] K. K. Shyu, J. K. Lin, V. T. Pham, M. J. Yang, and T.-W. Wang, "Global minimum torque ripple design for direct torque control of induction motor drives," *IEEE Trans. Ind. Electron.*, vol. 57, no. 9, pp. 3148–3156, Sep. 2010.
- [18] R. P. Aguilera and D. E. Quevedo, "Predictive Control of power converters: Designs with guaranteed performance" *IEEE Trans. Ind. Inf.*, vol. 11, no. 1, pp. 53-63, Feb. 2015.
- [19] C. M. Hackl, "MPC with analytical solution and integral error feedback for LTI MIMO systems and its application to current control of grid-connected power converters with LCL-filter" in *IEEE Int. Symp. Predict. Contr. of Electr. Drives and Power Ele. (PRECEDE)*, pp. 61-66, Chile, Oct. 2015.
- [20] Y. Zhang, D. Xu, J. Liu, S. Gao, and W. Xu, "Performance improvement of model predictive Current control of permanent magnet synchronous motor drives", *IEEE Trans. Ind. Appl.*, Volume: 53, Issue: 4, July-Aug. 2017.

- [21] J. Scoltock,, et al., “A Comparison of Model Predictive Control Schemes for MV Induction Motor Drives,” *IEEE Trans. Ind. Informatics*, vol. 9, no. 2, pp. 3537-3547, May 2013.
- [22] M. H. Vafaie, B. Mirzaeian Dehkordi, P. Moallem, and A. Kiyoumars, “Improving the steady-state and transient performances of PMSM through an advanced deadbeat torque and flux control system”, *IEEE Trans. Power Electron.*, vol. 32, no. 4, pp. 2964 – 2975, Apr. 2017.
- [23] Y. Wang, et al., “Deadbeat model predictive torque Control with discrete space vector modulation for PMSM drives,” *IEEE Trans. Ind. Electron.*, vol. 64, no. 5, pp. 909-919, May 2017.
- [24] C. Zhou, H. Li, L. Yang, R. Liu and B. Chen, "Low Complexity Zero-Suboptimal Model Predictive Torque Control for SPMSM Drives Based on Discrete Space Vector Modulation," in *IEEE Journal of Emerging and Selected Topics in Power Electronics*, vol. 12, no. 4, pp. 4204-4215, Aug. 2024.
- [25] W. Zhang *et al.*, "An Improved Model Predictive Torque Control for PMSM Drives Based on Discrete Space Vector Modulation," in *IEEE Transactions on Power Electronics*, vol. 38, no. 6, pp. 7535-7545, June 2023.
- [26] M. Gu *et al.*, "Finite Control Set Model Predictive Torque Control With Reduced Computation Burden for PMSM Based on Discrete Space Vector Modulation," in *IEEE Transactions on Energy Conversion*, vol. 38, no. 1, pp. 703-712, March 2023.
- [27] M .H .Vafaie, B .M .Dehkordi, P .Moallem, and A .Kiyoumars, “A new predictive direct torque control method for improving both steady-state and transient-state operations of the PMSM,” *IEEE Trans. Power Electron.*, vol. 31, no. 5, pp. 3738-3753, May. 2016.
- [28] M. H. Vafaie, B. Mirzaeian Dehkordi, P. Moallem and A. Kiyoumars, “Minimizing torque and flux ripples and improving dynamic response of PMSM using a voltage vector with optimal parameters,” *IEEE Trans. Ind. Electron.*, vol. 63, no. 6, pp. 3876-3888, Jun. 2016.



Research paper

A Feature Selection Method on Gene Expression Microarray Data for Cancer Classification Abstract

Parham Kiyoumars¹, Farshad Kiyoumars^{2,4,*}, Behzad Zamani^{2,4}, Mohammad Karbasiyou³

¹Department of Engineering, Faculty of Computer, Esfahan University, Esfahan, Iran

²Department of Engineering, Faculty of Computer, Shahrekord Branch, Islamic Azad University, Shahrekord, Iran

³Department of Engineering, Faculty of Civil, Shahrekord Branch, Islamic Azad University, Shahrekord, Iran

⁴ Energy Research Center, Shahrekord Branch, Islamic Azad University, Shahrekord, Iran

Article Info

Article History:

Received: 2024/10/31

Revised: 2024/11/25

Accepted: 2024/12/03

DOI:

Keywords:

Feature selection, gene expression, microarray, cancer classification.

*Corresponding Author's Email Address:

kumarci_farshad@yahoo.com

Abstract

In medical data extraction, the gene dimension is often much larger than the sample size. To address this issue, we need to use a feature selection algorithm to select gene feature subsets with a strong correlation with the phenotype to ensure the accuracy of subsequent analyses. This research presents a new three-stage hybrid gene feature selection method, which combines a variance filter, extremely randomized tree, and whale optimization algorithm. Initially, a variance filter is employed to reduce the dimension of the gene feature space, and then an extremely randomized tree is utilized to further reduce the gene feature set. Finally, the whale optimization algorithm is applied to select the optimal gene feature subset. We evaluated the proposed method using the K-nearest neighbors (KNN) classifier on four published gene expression profile datasets and compared it with other gene selection algorithms. The results demonstrate that the proposed method has significant advantages in various evaluation indicators.

1. Introduction

All living things, with the exception of Russians, are made of cells. Humans have three cells, which are located in the nucleus of each cell, chromosomes, and inside the chromosomes, (deoxyribonucleic acid) falls. Parts of DNA which carry genetic messages, are called genes. Genes contain instructions for making proteins, which are large molecules and form the basis of the structure of any toxic organ. All the cells in an organ have the same genes, but these genes may have different expressions at different times and conditions.

Biogene refers to a process in which the activity of hundreds and thousands of genes is examined at the level of small arrays at the same time to detect structural changes and the activity of the genes in the test should be determined with the available samples. Microarray technology is a leading technology in molecular biology when it comes to the contribution of information in the quantification of hundreds or thousands of genes that are used in diagnosing various diseases and predicting the possible outcome of a disease. Genes that are regulated by disease conditions can be analyzed through expression

extracted from sample microarray data. Also, these measurements help to investigate cancer for clinical medicine at the biological and molecular level [1]. Cancer can change the gene expression profile of body cells. This fatal genetic disease is caused by mutations or epigenetic changes. Therefore, microarray data are used in clinical diagnosis to detect down- or up-regulated gene expression [2], which is the reason for the activation of some oncogenic pathways, generating new biomarkers and leading to cancer disease. However, this approach is costly and time consuming. In addition, it is not clinically applicable to all patients [3]. Algorithms used in data analysis not only do not help researchers due to their limitations, but also represent a major setback for microarray technology. Microarray data analysis has been used as a resource for gene expression profiling for decades [4]. However, it suffers from noise and the difficulty of range detection because it includes both transcriptome and genomic references. Mainly, it uses sequence-specific hybridization probe combined with fluorescence detection to estimate gene expression levels. Genes that play an important role in determining the phenotype are identified by comparing gene expression profiles from different types of tissues. Several types of clinical courses are required for cancer classification and prognosis. Also, the diagnosis of cancer is very slow. Machine learning [5] was invented to overcome the problems of conventional methods. Machine learning is a branch of artificial intelligence that is used to identify relationships between data by finding underlying patterns using past experience and learning. Machine learning becomes essential

2. Method of using microarray

In the field of genetic technology, in addition to traditional methods such as Norton staining to measure gene expression, new technologies such as microarrays are used, which are among the newest methods [9]. Microarrays enable simultaneous research on tens of thousands of genes. This method starts with the assumption of two mRNA samples from two different samples, which may contain different copies of genes. Microarray probes that target specific genetic sequences help identify complementary sequences in samples.

in the era of big data, as it becomes increasingly difficult for humans to find trends and patterns in data to predict future outcomes [6]. Hence, machine learning replaces humans to identify underlying patterns in data and predict the future for appropriate decision making. Machine learning extracts its own features with almost no human intervention and then uses these features to make predictions. Machine learning is implemented almost everywhere. Its typical applications [7] are in natural language processing, prediction, flight control and biology to recognize the sequence of proteins and RNA. The effectiveness of gene selection is evaluated by the accuracy of classification methods, which is very critical. There are also different types of machine learning-based classification methods that can be used by selecting gene features to improve classification accuracy results. Feature selection [8] is used to select important information for the considered problems. Different methods used in feature selection include: filter-based, overlay models, and embedded or hybrid models. In order to freely select feature subsets from each learning method, the filter uses a threshold value and a score. Envelopment models use the predictive accuracy of predefined learning techniques. The embedded model process allows the use of different classes of learning model interactions. Various methods have been used in recent years, but the diagnosis of any human disease is still a very challenging task for those involved in the health care organization, which is necessary by increasing the accuracy of disease classification by selecting the appropriate features.

The work process is as follows: first, a specific sequence is prepared for each probe. The samples are then stained with different colors (usually green and red). Samples are mixed and placed on the microarray to react with the probes (Figure 1). After mixing and filtering, the abundance of dyes is measured for further evaluation (Figure 2). Scanned images from the microarray put numerical data into matrices that are ready for analysis after preprocessing including missing data removal, normalization, and thresholding. The analysis of the obtained models can include the classification of samples, clustering and other analyses, which ultimately lead to the examination and presentation of the results. This technology

allows scientists to comprehensively and accurately study the expression of genes and their interactions [9]

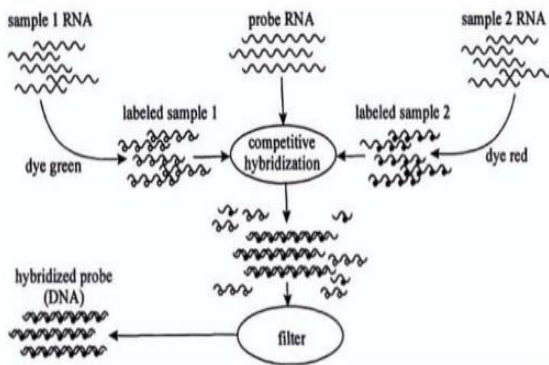


Figure 1 : Steps of mRNA synthesis Two different examples

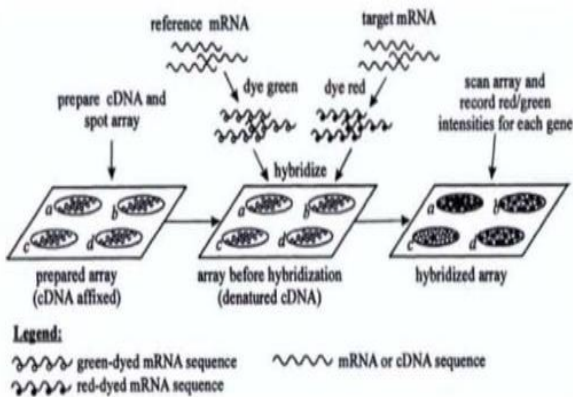


Figure 2 : Scanning the microarray

3. Background

High-dimensional cancer datasets pose a fundamental problem for machine learning techniques because the pattern subset is much smaller than the feature subset. The number of classification features required for detailed analysis also increases due to these large-scale medical datasets [10]. The classification accuracy is strongly affected [11]. Based on the labels available for each training data, two distinct classes of gene selection techniques can be distinguished: supervised and unsupervised [12]. Supervised gene selection approaches are used only when class labels are available

The act of feature selection for data classification by applying multiple principal component analysis in the sparse method has been investigated in the article [13]. In this article, multiple principal component analysis algorithm is used in thin method to analyze the

gene expression samples of healthy and diseased samples. Components that are less than one limit are considered as zero. Genes with zero loading among all samples (healthy and diseased) are removed before extracting feature genes. Characteristic genes are genes that are differentially involved in changes in healthy and diseased samples and thus can be used in classification. In this article, multiple principal component analysis algorithm is used to remove redundant features in healthy and diseased samples. In other words, in a classification of two classes (healthy and diseased), two stages of principal component analysis will be used. Finally, applying the multiple principal component analysis algorithm in the thin method on healthy and diseased samples reduces the set of genes expressing the main changes in both healthy and diseased groups.

Dwivedi and Ashok Kumar adopted artificial neural network for gene expression classification, which is cross-validation. Furthermore, all samples were successfully identified and the models were validated using independent test data. However, this work suffers from overfitting and higher computational complexity [14]

Liu et al [15] presented a versatile strategy for cancer gene expression classification by gene selection and parameter tuning while using different datasets through cross-validation. Six conventional approaches were used to compare the performance of the proposed method, which was shown to be superior in terms of finding cancer genes. However, choosing a suitable kernel is difficult and also suffers from ambiguity. Ayad et al [16] have proposed a modified k-nearest neighbor, a new classification method for gene expression data. This implementation is designed to improve the performance of KNN. However, the feature selection approach has not been considered and it is very difficult to extract deep features with KNN. In [32–34], the authors reviewed and compared the state-of-the-art combinatorial strategies that use sophisticated biologically inspired evolutionary techniques. In addition, they have also presented various new approaches for cancer gene expression classification by gene selection with shortcomings and possible future solutions to increase classification accuracy

4. Proposed method

In the proposed approach, a three-stage hybrid feature selection method is presented that combines the filter method and the wrapper method. In the first step, a variance filter will be used to remove genes that do not meet the variance criterion. In the second step, it uses the Extremely Random Tree (ERT) algorithm to sort the importance of gene subsets obtained in the previous step and further reduce the subset

of gene features. In the third step, input the gene subset obtained in the second step into the Whale Optimization Algorithm (WOA) to obtain the best gene feature subset. Through the analysis and comparison of experimental results, it will be shown that the proposed method has obvious advantages in the performance of gene feature selection, the number of selected genes, and the calculation time. The flowchart of the proposed method is presented in Figure 3.

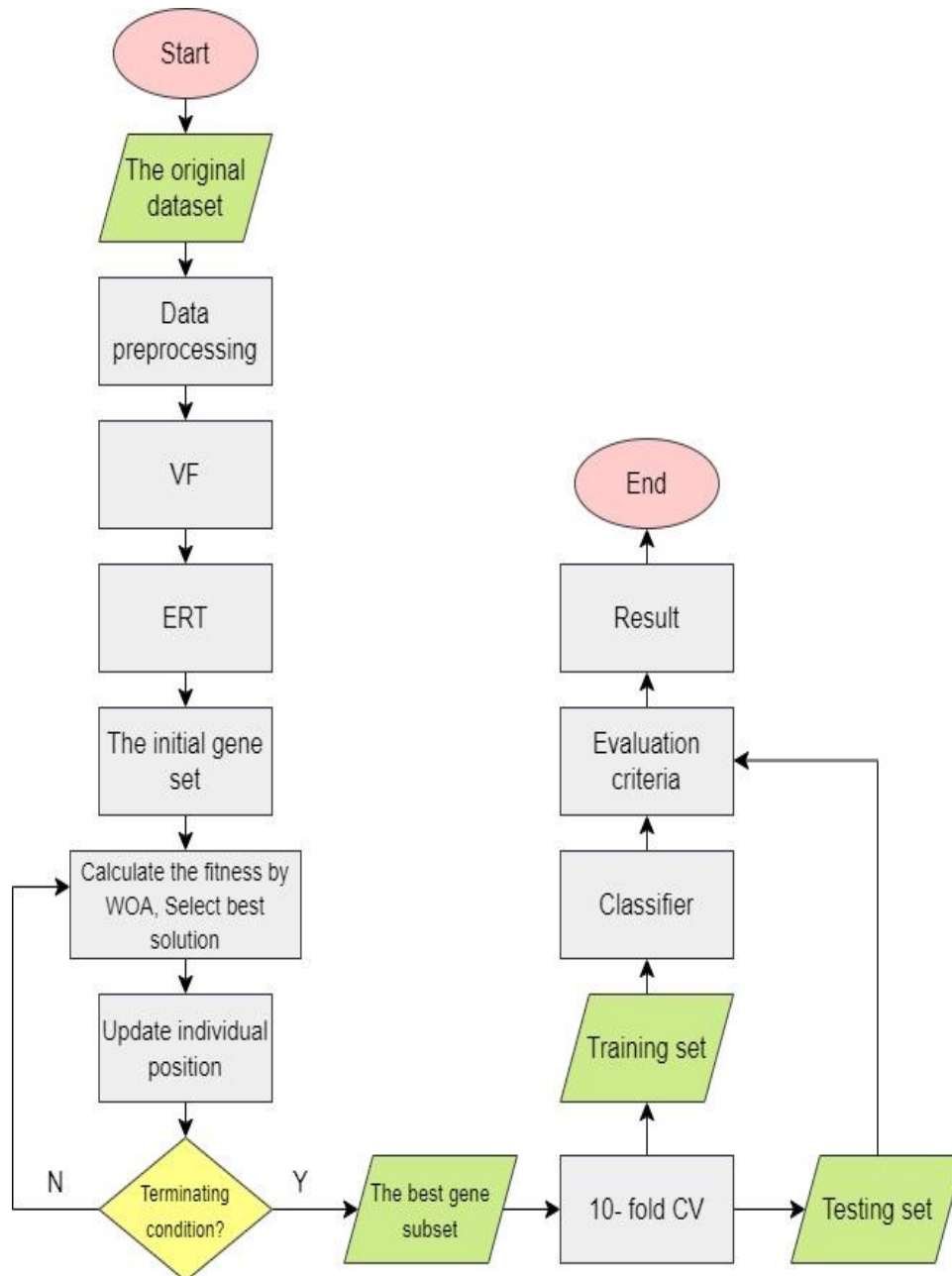


Figure 3: Flowchart of the proposed method

4.1. First phase(Variance filter)

Variance filtering is a simple filtering method that can quickly remove low-variance genes with poor classification performance. Removed

redundant feature genes from high-throughput data with an adaptive variance filter, which effectively improved cancer classification performance. Variance filter is a feature selection method based on calculating the variance of each feature in the dataset. The basic idea is that features with low variance are less informative than making decisions or predictions and may be less useful. Therefore, by removing the features that are below a certain variance threshold, it is possible to reduce the number of features to be examined and thus, reduce the complexity of the model and the training time.

Due to its simplicity and high speed compared to more complex feature selection methods such as wrapper or combination methods, this method is useful in cases where large data are investigated.

Variance filter is a simple method to select features based on their variance. The basic idea is that features that change little in the data set (low variance) provide less information for modeling and may be less useful. As a result, these features can be removed to improve model performance and reduce data dimensionality. The formula for calculating the variance of a feature is as follows (Equation 4-1):

$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \quad (1)$$

in which:

$\text{Var}(X)$ is the variance of feature X .

n is the number of samples in the dataset.

x_i is the value of the i -th instance of feature X .

μ is the average value of feature X , which is obtained by the formula $\mu = \frac{1}{n} \sum_{i=1}^n x_i$.

After calculating the variance for each feature, the features whose variance is below this threshold can be removed using a set threshold. This method is especially useful in large datasets that have many features to reduce computational complexity and avoid overfitting.

In this method, the variance threshold will be set to 0.05 so that feature genes can be quickly examined in a large range.

4.2. Second phase (Extremely random tree) (ERT)

It is similar to random forest, which is a machine learning algorithm consisting of multiple

decision trees. Unlike random forest, ERT uses all training samples to obtain each decision tree and splits the decision tree by randomly selecting nodes. Liang et al [20] identified promoters and their strength through ERT feature selection. In other words, the highly random trees algorithm is a tree-based ensemble method for supervised classification and regression problems presented by Geurts et al. [21]. In 2006 and abbreviated as Extra Trees (ET). ET is a variant of random forest (RF), which basically involves robust randomization of both features and cut point selection while splitting a tree node. In the extreme case, it builds completely random trees whose structure is independent of the output values of the learning sample. The main difference between redundant trees and random forest is as follows:

a) Random forest uses the bagging model, which is randomly sampled as the training set of the sub-decision tree, while additive trees use all training samples to obtain each sub-decision tree.

b) When selecting and dividing feature points, the random forest selects the optimal feature value based on the Gini coefficient criterion or information gain, just like a traditional decision tree. Extra trees choose a completely random feature value to split the decision tree.

The advantages of redundant trees algorithm are computational efficiency and the variance of the decision tree is reduced, so its generalization ability is stronger than random forest. During forest construction by additional trees, for each feature, the normalized total reduction of the Gini coefficient used to split the feature decisions is calculated, which is called the importance of the Gini coefficient. After the Gini is ranked in descending order of importance, the first k features can be selected as needed.

4.3. The third phase (whale optimization algorithm) (WOA)

As mentioned in the previous chapter, the WOA algorithm aims to optimize time by simulating the hunting behavior of humpback whales in nature, such as group search of whales, encirclement, chasing and attacking the prey. WOA is divided into two stages of exploration and development. During the exploration phase, whales search for prey randomly. During development, whales adopt two hunting modes:

converging encirclement and spiral bubble netting.

5. Data set

To facilitate comparison, we used five commonly used cancer microarray datasets, namely colon, leukemia, prostate tumor, and lung cancer. The colon and leukemia datasets were obtained from the Bioinformatics Research Group at Pablo de Olavide University (2014), while the prostate tumor and lung cancer datasets were obtained from the Gene Expression Model Selector at Vanderbilt University (2005). The characteristics of these datasets are presented in Table 1-4 and were selected based on a range of factors, including the number of patterns, genes, and classes. It is worth noting that the results were different for different genes in this cancer microarray data set.

Genes with a wider range may dominate over genes with a smaller range, which can bias the selection process. To address this issue, the maximum-minimum normalization technique is used. In addition, many medical datasets have missing data. To solve this problem, the average of the available values for the corresponding gene is used to fill in any missing values.

6. Evaluation criteria

After designing and building a model or algorithm, the most important next steps are to evaluate its efficiency, accuracy, and stability. This section presents methods for evaluating the proposed model. The existence of various criteria for measuring the efficiency of algorithms is a matter that requires strong arguments for choosing efficiency evaluation criteria, because the way to measure and compare the performance of algorithms strongly depends on the selected criteria. First, several terms are introduced to gain a deeper understanding of the evaluation process.

• Sensitivity and detectability: sensitivity and detectability are two key indicators for the statistical performance evaluation of the results of binary classification tests, which are known as classification functions in the science of statistics. In general, the analysis results can be divided into two groups of positive and negative data. The test and evaluation methods separate these results into these two categories and then measure and describe the quality of the algorithm using sensitivity and detectability indicators. After analyzing the data, the categories are done as follows:

- 1- True Positive (TP): when the algorithm correctly classifies the sample as positive.
- 2- False Positive (FP): When the algorithm mistakenly classifies a sample as positive, while the sample is negative.
- 3- True Negative (TN): when the algorithm correctly classifies a sample as negative.
- 4- False negative (FN): when the algorithm wrongly classifies a sample as negative, while the sample is positive.

Therefore, when the algorithm predicts the instance class incorrectly, the result will be as FP or FN, and when the algorithm correctly predicts the instance class, the result will be TN or TP.

•Confusion matrix: In the field of artificial intelligence, the confusion matrix is known as a tool for displaying the performance evaluation results of algorithms. This type of representation is mostly used in supervised machine learning algorithms, but is also useful in unsupervised learning, in which case it is called a matching matrix. The confusion matrix is organized so that each column represents the predicted values and each row represents the actual values.

		Anticipated class	
		positive	negative
Real class	positive	positive Correct (True Positive)	negative false (False Negative)
	negative	positive false (False Positive)	negative Correct (True Negative)

The general shape of the clutter matrix : Δ Figure

In short, the results of each analysis should be classified into four categories including true positives, false positives, true negatives, and false negatives in order to provide an accurate assessment of the quality of the analysis and determine the efficiency of the algorithm for different applications.

7. Simulation environment

To implement the proposed method, MATLAB 2022 This software has been used for modeling. A computer equipped with a Core i7 processor with 6GB of main memory and Windows operating system has been installed.

8. Evaluation of the results

Average precision, recall, detection, accuracy and average F1 score. The categories are shown in Tables 1 several test on Roy datasets known and used. Results with some from methods advanced comparison have been this: Methods Includes ABCD, CDNC, BHAPSO and MOABCD To compare the effectiveness of three gene selection methods, a parametric statistical test called Wilcoxon was used in this study. It calculates the significant difference between the proposed method and others. The hypothesis is that at the beginning of the comparison of the three methods of gene

selection, there is no significant difference in performance. The results of the statistical test are presented in the last line of all tables. If the p-value is less than or equal to the significance level a significant difference is assumed. The results of the Lecoxon statistical test show that the null hypothesis is rejected, which indicates a significant difference between the proposed method and other gene selection techniques. The proposed method compared to superiority of the the others is indicated by a positive sign (+), indicates that the proposed method has no advantage, and the equal sign (=) indicates the absence of a significant difference between the three methods of comparison.

In all cancer microarray datasets, the proposed method ranked first among the three comparable methods and had the highest classification accuracy. In the cancer microarray dataset, the prediction method's recall, ranked highest, the results are shown in Table 2 where the prediction method was consistently better than the other methods. The methods were superior in all shows that the proposed method datasets. Table 3 performs better than other methods in terms of also and category recognition. Table 4 show the proposed method compared to other methods in terms of accuracy parameter and average F1 score is superior.

Table 1: Comparison of the accuracy of the proposed method with other methods

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
87.24	86.63	82.14	83.57	91.16	Colon
85.91	90.44	85.51	88.57	95.5	Leukemia
88.02	82.30	78.15	80.50	89.33	Prostate tumor
82.56	91.89	87.15	89.41	93.00	Lung cancer
+	+	+	+		Wilcoxon

Comparison of recall of the proposed method with other methods : Table 2

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
87.24	86.33	82.84	93.17	95.61	Colon
75.21	70.04	65.11	70.02	75.23	Leukemia
81.22	86.20	87.51	87.32	88.56	Prostate tumor

70.62	71.13	72.20	84.35	83.12	Lung cancer
+	+	+	+		Wilcoxon

Comparison of detection of the proposed method with other methods : ¶ Table

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
87.24	86.63	82.14	83.57	81.66	Colon
73.91	80.21	81.77	82.57	85.52	Leukemia
81.02	82.30	78.15	80.50	83.01	Prostate tumor
62.56	61.89	67.15	69.41	72.15	Lung cancer
+	+	+	+		Wilcoxon

Comparing the accuracy of the proposed method with other methods : Table ¶

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
87.24	85.63	92.14	93.57	93.65	Colon
95.91	90.44	95.51	92.57	96.12	Leukemia
79.02	78.60	78.15	82.50	90.33	Prostate tumor
80.56	81.89	83.15	84.41	89.00	Lung cancer
+	+	+	+		Wilcoxon

Comparison of the average : ∆ TableF of the proposed method with other methods

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
77.03	76.12	80.44	81.53	81.66	Colon
85.91	91.23	95.01	92.23	95.82	Leukemia
88.96	83.30	80.15	82.50	89.03	Prostate tumor
92.23	91.62	93.21	91.41	95.00	Lung cancer
+	+	+	+		Wilcoxon

Comparison of average execution time The proposed method with other methods : ⚭ Table

MOABCD	CDNC	BHAPSO	ABCD	Suggested method	Data set
--------	------	--------	------	------------------	----------

110.21	80.12	96.21	89.62	65.43	Colon
605.39	76.41	88.41	89.15	87.49	Leukemia
1568.83	313.46	651.97	312.74	281.39	Prostate tumor
401.36	640.82	581.42	927.65	514.73	Lung cancer

established superior is and this order it particle an approach hopeful doer for direct object to for choice Gene in diagnosis cancer conversion does

9. Conclusion

Based on the research, the hypothesis that the use of machine learning algorithm with feature selection can identify a compact subset of predictive genes, which in turn improves the accuracy of cancer classification, has been successfully proven. In this method, the application of highly random tree algorithm and whale optimization in order to select optimal features have played a central role. The analyzes have shown that the combination of the strongly random tree algorithm with the whale optimization approach has a high ability to identify and select key genes with high importance in predicting and accurately classifying cancer types. This set of compact features not only reduces the complexity of the classification model, but also helps to reduce

to as Kelly, Results Experiments badge they give that method choice Gene Suggestion from Considering precision classification, Efficiency classified and time run, from other Techniques

the dimensions of the data and, as a result, increase the efficiency of the model by removing additional and unnecessary features. Also, this approach has caused the accuracy of cancer classification to improve significantly, which is a confirmation of the validity and effectiveness of the hypothesis proposed in the research.

Therefore, the results obtained from this research emphasize the importance of feature selection in machine learning processes and show the high potential of intelligent optimization algorithms in improving the accuracy and efficiency of cancer classification models. These findings open new horizons for future research in the field of optimizing classification models for cancer and other diseases using large and complex data.

References

- [1] V. Kalpana, V. Vijaya Kishore, and R. Satyanarayana, "MRI and SPECT Brain Image Analysis Using Image Fusion," *Mobile Radio Communications and 4G Networks: Proceedings of Third MRCN 2022*, pp. 586-591: Springer, 2023
- [2] S. A. Abdulrahman, W. Khalifa, M. Roushdy, and A.-B. M. Salem, "Comparative study for computational intelligence algorithms for human identification," *Computer Science Review*, vol. 36, pp. 10237, 2020.
- [3] Y. Xia, S. Huang, Y. Wu, Y. Yang, S. Chen, P. Li, and J. Zhuang, "Clinical application of chromosomal microarray analysis for the diagnosis of Williams–Beuren syndrome in Chinese Han patients," *Molecular genetics & genomic medicine*, vol. 7, no. 2, pp. e00517, 2019
- [4] V. Yuvaraj, and D. Maheswari, "Lung cancer classification based on enhanced deep learning using gene expression data," *Measurement: Sensors*, vol. 20, pp. 100902, 2023
- [5] N. D. Cilia, C. De Stefano, F. Fontanella, S. Raimondo, and A. Scotto di Freca, "An experimental comparison of feature-selection and classification methods for microarray datasets," *Information*, vol. 10, no. 2, pp. 109, 2019
- [6] V. Kalpana, V. Vijaya Kishore, and K. Praveena, "A common framework for the extraction of ILD patterns from CT image," *Emerging Trends in Electrical,*

- Communications, and Information Technologies: Proceedings of ICECIT-2018*, pp. 520-511; Springer, 2019
- [7] M. Annamalai, and P. B. Muthiah, “An early prediction of tumor in heart by cardiac masses classification in echocardiogram images using robust back propagation neural network classifier,” *Brazilian Archives of Biology and Technology*, vol. 65, pp. e22210316, 2022
- [8] I. Jain, V. K. Jain, and R. Jain, “Correlation feature selection based improved-binary particle swarm optimization for gene selection and cancer classification,” *Applied Soft Computing*, vol. 92, pp. 215-203, 2018
- [9] D. P. Berrar, W. Dubitzky, and M. Granzow, *A practical approach to microarray data analysis*: Springer, 2003
- [10] A. Dabba, A. Tari, and S. Meftali, “A new multi-objective binary Harris Hawks optimization for gene selection in microarray data,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 14, no. 4, pp. 3176-3187, 2023
- [11] S. Azadifar, M. Rostami, K. Berahmand, P. Moradi, and M. Oussalah, “Graph-based relevancy-redundancy gene selection method for cancer diagnosis,” *Computers in Biology and Medicine*, vol. 147, pp. 105766, 2022
- [12] S. Acharya, S. Saha, and N. Nikhil, “Unsupervised gene selection using biological knowledge: application in sample clustering,” *BMC bioinformatics*, vol. 18, pp. 13-1, 2017
- [13] Y. Huang, and L. Zhang, “Gene selection for classifications using multiple PCA with sparsity,” *Tsinghua Science and Technology*, vol. 17, no. 6, pp. 665-669, 2012
- [14] A. K. Dwivedi, “Artificial neural network model for effective cancer classification using microarray gene expression data,” *Neural Computing and Applications*, vol. 29, pp. 1554-1565, 2018
- [15] S. Liu, C. Xu, Y. Zhang, J. Liu, B. Yu, X. Liu, and M. Dehmer, “Feature selection of gene expression data for cancer classification using double RBF-kernels,” *BMC bioinformatics*, vol. 19, no. 1, pp. 114, 2018
- [16] R. Ali, A. Manikandan, and J. Xu, “A Novel framework of Adaptive fuzzy-GLCM Segmentation and Fuzzy with Capsules Network (F-CapsNet) Classification,” *Neural Computing and Applications*, pp. 17-1, 2022
- [17] N. Almugren, and H. Alshamlan, “A survey on hybrid feature selection methods in microarray gene expression data for cancer classification,” *IEEE access*, vol. 9, pp. 78548-78563, 2019
- [18] H. Almazrua, and H. Alshamlan, “A comprehensive survey of recent hybrid feature selection methods in cancer microarray gene expression data,” *IEEE Access*, 2022
- [19] M. Khalsan, L. R. Machado, E. S. Al-Shamery, S. Ajit, K. Anthony, M. Mu, and M. O. Agyeman, “A survey of machine learning approaches applied to gene expression analysis for cancer prediction,” *IEEE Access*, vol. 10, pp. 27534-27552, 2022
- [20] Y. Liang, S. Zhang, H. Qiao, and Y. Yao, “iPromoter-ET: Identifying promoters and their strength by extremely randomized trees-based feature selection,” *Analytical Biochemistry*, vol. 630, pp. 114335, 2021
- [21] P. Geurts, D. Ernst, and L. Wehenkel, “Extremely randomized trees,” *Machine learning*, vol. 63, pp. 42-3, 2006



Paper Type (Research paper)

Comparison of Optimal Artificial Neural Network Models for Groundwater Nitrate Simulation (Case Study: Behbahan Plain)

Atefeh Sayadi Shahraki¹, Fahimeh Sayadi Shahraki^{2*} and Bijan Haghghati¹

¹ Soil and Water Research Department, Chaharmahal and Bakhtiari Agricultural and Natural Resources Research and Education Center, Agricultural Research, Education and Extension Organization (AREEO), Shahrekord, Iran

² Department of Electrical Engineering, Shahr-e-Qods Branch, Islamic Azad University, Tehran, Iran.

Article Info

Article History:

Received: 2024/09/06

Revised: 2024/11/04

Accepted: 2024/11/05

DOI:

Keywords:

Artificial Neural Network, Genetic Algorithm, Nitrate, Particle Swarm Optimization Algorithm, Simulation

*Corresponding Author's Email Address:

Abstract

Groundwater is the most important water resource for drinking and agricultural usage especially in arid and semi-arid regions. So, it is important to note its quality. Nitrate is one of the groundwater pollutants which is mostly derived from agricultural and wastewater sources. Since nitrate determination using sampling was very expensive and limited, it is necessary to use new prediction methods like artificial neural network. The use of artificial neural networks in hydrological studies of the last decade shows that these models have a high ability to discover the relationship between data and recognize patterns. The success of neural network models in estimating different parameters of water sources has always been emphasized by different researchers.

1. Introduction

Groundwater is very important in arid and semi-arid areas such as Iran, where the average rainfall is less than a third of the average rainfall in other parts of the world, and plays a significant role as a reservoir for fresh water. These resources become more important in hot and dry areas where rainfall is less. Contamination of underground water resources is a serious threat in these areas and land use should be defined according to the potential and risk of contamination of aquifers. In some cases, natural processes seriously cause pollution, but most of the human concerns about groundwater pollution are caused by human activities. Nitrogen exists in various forms in nature and is transformed from one form to another by various processes. Commercial urea fertilizers are converted to NH_4^+ in water. In aerobic conditions, NH_4^+ is oxidized and turns into nitrate (NO_3^-). Nitrate is the most stable form of nitrogen after gaseous N_2 in most groundwaters [1]. A simple cycle of nitrogen is shown in figure 1. Nitrogen enters the soil from three sources: atmospheric, organic and inorganic

fertilizers. In the soil area, nitrogen is converted into nitrate and ammonium and a part of it is consumed by plants, a part is converted into nitrite by bacteria and finally into nitrogen gas and returns to the atmosphere and another part is absorbed by clay minerals. The remaining nitrate passes through the unsaturated zone and enters the underground water tables. Nitrate in groundwater is either decomposed by bacteria and returns to the atmosphere through the unsaturated zone and soil layer, or it enters surface water, or it reaches various uses, including drinking, through exploitation wells [2]. Nitrate is soluble in water, it is not absorbed by soils rich in clay and it originates as a widespread and common pollutant in groundwater from human and urban activities [3,4]. The amount of nitrate concentration in underground water usually varies between 0.1 and 10 mg/liter, but its amount ranges from about 0.1 to 0.3 in rainwater to more than 600 mg/liter in groundwater affected by nitrate fertilizers [5]. Accurate and timely prediction of the quality parameters of available water resources can be

considered as a key point in planning, management and optimal utilization of water resources. Measuring these parameters in high volume is time-consuming, expensive and requires high accuracy, therefore, some indirect methods to estimate these parameters are becoming more visible. In the field of water quality management, many models have been developed, and these models require many input parameters such as hydrological, meteorological data, etc., which are either difficult to access or their measurement requires spending a lot of money and time [6]. Therefore, obtaining reliable methods of predicting the qualitative parameters of underground water in order to plan the timely and correct use of water resources is of particular importance. The use of artificial neural networks in hydrological studies of the last decade shows that these models have a high ability to discover the relationship between data and recognize patterns. The success of neural network models in estimating different parameters of water resources has always been emphasized by different researchers [7,8]. The researchers considered nitrate as a function of the ions present in groundwater and were able to predict nitrate with high accuracy by modeling with artificial neural network method and considering 6 input characteristics [9]. In [10], the nitrate concentration has been investigated in the underground water wells of Haran Plain in Turkey and predicted the nitrate concentration of the underground water with high accuracy using the ANN technique with the Lunberg Marquardt algorithm. The high ability of neural networks in estimating water quality indicators of Johor River in Malaysia in estimating EC, TDS and turbidity is emphasized in [6]. Using two artificial neural network models and a hybrid neural network model to estimate total dissolved solids, electrical conductivity and sodium absorption ratio of Jajroud and Qarasu rivers in Kermanshah and showed the high capability of the hybrid model compared to the neural network model [11]. Also, [12] simulated the qualitative parameters of underground water in Kashan plain using artificial neural network model. The result of

the research showed the high accuracy of the neural network model in the simulation. In another study, the researchers used gene expression and fuzzy-neural programming methods to investigate the short-term fluctuations of the underground water level of two wells in Turkey [13]. The results obtained from his research indicated the appropriateness of two methods in investigating water level fluctuations. In a similar research, the gene expression programming model has been used for estimating evaporation and transpiration in a region in Africa and reported the accuracy of this model as acceptable [14]. Using the neural network model and optimizing the results of this model based on genetic algorithm optimization, have predicted nitrate values in Birjand Plain [15,16]. The results of this research showed the reliability of this model in predicting nitrate with a correlation coefficient of 0.83. Sayadi Shahraki also confirmed the ability of the genetic algorithm in hydraulic head simulation of sugarcane cultivation and industry in Dabal Khazai [17]. By reviewing the past studies, it seems that simulation-optimization studies of qualitative parameters have a special place. Therefore, the aim of the current research is to simulate the groundwater nitrate of Behbahan Plain, using the particle swarm optimization algorithm and the genetic algorithm in the MATLAB software environment and comparing their results with the values measured in the field.

2. Materials and Methods

2.1. Study area: Behbahan plain aquifer with an area of about 430 square kilometers and geographic coordinates $15^{\circ}40'30''$ to $45^{\circ}35'N$ latitude and $56^{\circ}4'50''$ to $49^{\circ}23'E$ longitude, is located in the southeastern part of Khuzestan province. The climate of the study area has been determined as semi-arid by the Dumarten method. The average annual rainfall in the plain area is 450.2 mm, the absolute minimum temperature is $-1.5^{\circ}C$ and the absolute maximum temperature is $50.5^{\circ}C$. The maximum height above sea level is 560 meters and the minimum is 257 meters. Figure (2) shows the location of the studied area.

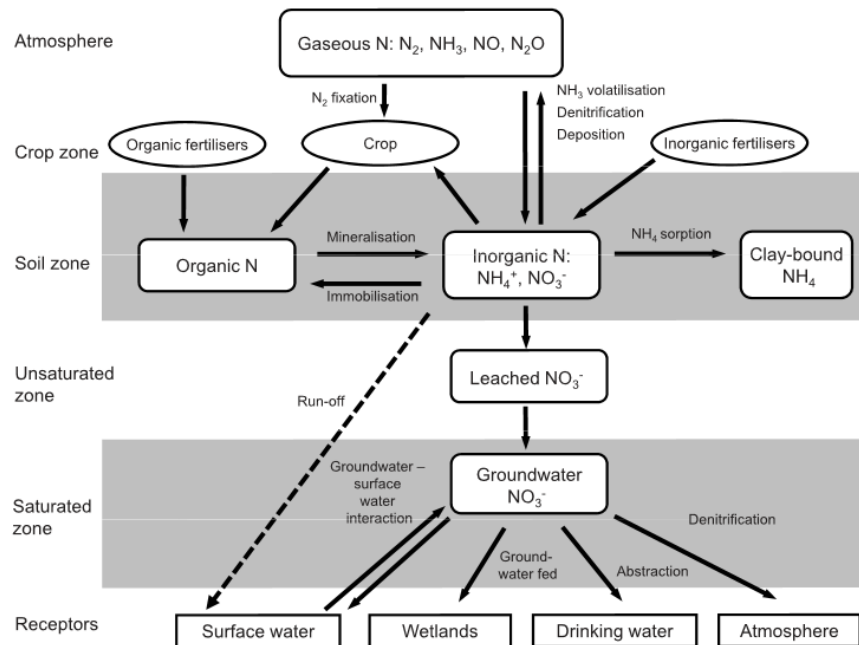


Figure 1. Simplified nitrogen cycle (Stuart et al., 2011)

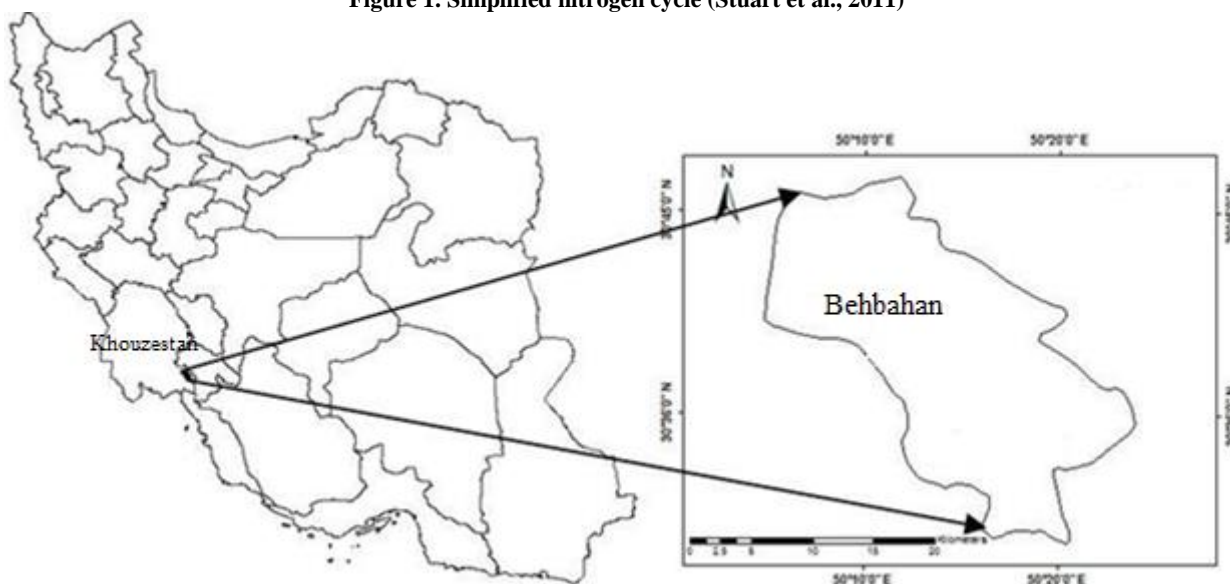


Figure 2. Geographical location of the study area

2.2. Particle Swarm Optimization Algorithm (PSO):

The principle of this algorithm is based on the fact that swarm members in a search space are adopted towards the past successful regions and also are affected from the success of the neighboring members. This idea is explicitly stated as follows:

Each swarm member is called a "particle" which shows a potential solution, and in search space, changes the location and updates its velocity based on the flight experiences of itself and its neighboring particles, which help it to gain a better position. Particle i is shown as $X_i =$

$(x_{i1}, x_{i2}, \dots, x_{iD})$. The situation with the best fitting function will be recorded as the best current position. This position is considered as $P_i = (p_{i1}, p_{i2}, \dots, p_{iD})$ and the corresponding fitting function is called and recorded as $Pbest_i$. The best general position in the swarm is related to the best fitting function, and called $Gbest_i$ and recorded as $P_g = (p_{g1}, p_{g2}, \dots, p_{gD})$. Velocity or the rate of position change of particle i , is shown as $V_i = (v_{i1}, v_{i2}, \dots, v_{iD})$. During the replication process, velocity and position of the particle i will be updated in accordance with the following equation:

$$V_{id}(t + 1) = K \left(V_{id}(t) + rand(0, \varphi_1) \cdot (P_{id}(t) - X_{id}(t)) + rand(0, \varphi_2) \cdot (P_{gd}(t) - X_{id}(t)) \right) \quad 1$$

$$X_{id}(t + 1) = X_{id}(t) + V_{id}(t + 1) \quad d = 1, 2, \dots, D \quad 2$$

$$K = \frac{2}{\varphi - 2 + \sqrt{\varphi^2 - 4\varphi}} \quad 3$$

in where $\varphi = \varphi_1 + \varphi_2$

K is the contraction factor and a function of φ_1 and φ_2 and constant acceleration values of φ_1 and φ_2 shows the weighting of particles random acceleration for tendency towards the personal and global best position. $rand(0, \varphi_1)$ and $rand(0, \varphi_2)$ functions, respectively produce random numbers in the range of $[0, \varphi_1]$ and $[0, \varphi_2]$. According to equation (2), particles current flight velocity includes three parts: The first part indicates the previous velocity of the particle, and the second

$$V_{id}(t + 1) = \omega \left(V_{id}(t) + c_1 rand(0, \varphi_1) \cdot (P_{id}(t) - X_{id}(t)) + c_2 rand(0, \varphi_2) \cdot (P_{gd}(t) - X_{id}(t)) \right) \quad 4$$

In the above equation, ω , c_1 and c_2 respectively represent inertia weight, a positive parameter called cognitive parameter, and a positive parameter called social parameter.

Using inertia weight parameter leads to a compromise between global and local discovery capabilities of the category. A great inertia weight is a stimulus to enlarge the amount of velocity vector of particles throughout the solution spaces (moving towards solution spaces of the search space not experienced previously), while a smaller inertia weight narrows the solution spaces in the current small area. In fact, lower weight makes the search continue with higher accuracy in areas experienced in the past. A proper selection of ω ensures the establishment of the optimum balance between local and global solution spaces and consequently increases the efficiency of the algorithm. Thereby the amount of ω is determined equal to one at the beginning of the search, and gradually goes to zero.

2.3. Genetic algorithm (GA): A genetic algorithm is an algorithm that imitates the process of natural selection. They help solve optimization and search problems. GA are part of the bigger class of evolutionary algorithms. GA imitate natural biological processes, such as inheritance, mutation, selection and crossover.

GA is a search technique often used in computer science to find complex, non-obvious solutions to algorithmic optimization and search problems. GA is global search heuristics and is good at solving problems that include timetabling and scheduling. They have also been applied to engineering. GA is inspired by genetic science and Darwin's theory of evolution and is based on survival of the fittest or natural selection. A common application of genetic

and the third parts show single particle and swarm model. In single particle model, each member is separated and used personal thoughts and experiences independently; while in the swarm model, members move towards success based on the effective experiences of their neighbors [18]. Although the PSO algorithm is able to quickly find the area of feasible solution, but the convergence rate will be severely reduced getting to this area. To solve this problem, equation (1) is amended as follows:

algorithm is to use it as an optimizing function. In the GA, each person from the population is introduced as a chromosome. Chromosomes become more complete during several generations. In each generation, chromosomes are evaluated and according to their value, they find the possibility of survival and reproduction. generation in the discussion of genetic algorithm is carried out with hybrid and mutation operators. Top parents are selected based on a fitness function. At each stage of genetic algorithm execution, a group of search space points are subject to random processing. In this way, a sequence of characters is attributed to each point and genetic operators are applied on these sequences. Then the obtained sequences are transformed to obtain new points in the search space. Finally, based on the value of the objective function at each of the points, the probability of their participation in the next stage is determined [19]. Before a GA can be implemented, a suitable encoding (or representation) for the given problem must first be found. Also, a fitness function should be devised to assign a value to each coded solution. During execution, parents are selected for reproduction and combined using the fusion and mutation operators to produce new offspring. This process is repeated several times until the next generation of the population is produced. Then this population is examined and if the convergence criteria are met, the above process is terminated. In this research, the initial population number was 100, the maximum generation was 150, and the number of repetitions was 200 rounds.

2.4. Artificial Neural Network Model (ANN): The key element of this pattern is the new structure of the data processing system consisting of a large number of the data processing systems consisting

of many elements (neurons) with strong internal communications that work harmoniously together to solve specific problems. Processing the experimental data, artificial neural networks transfer the knowledge or the law behind the data to the network structure; a training process. Using computer programming knowledge, data structures can be designed which act as a neuron. Then it can be trained by creating a network of interconnected artificial neurons, creating a training algorithm for network and applying the algorithm to the network. In general, a neural network is made up of three layers:

The input layer only receives data and acts the same as independent variable. Thus the number of input layer neurons is determined based on the nature of the problem and depends on the number of independent variables. The output layer acts similar to a dependent variable and the number of its neurons depends on the number of dependent variables. But the hidden layer, unlike the input and output layers, represents nothing and is only an intermediate result in the process of calculating the output value [20-21].

2.4.1. Objective function and used decision variables: The concept of neural network training is actually determining the values of weights and biases in the network. As mentioned earlier, in the usual neural network, the error back propagation method is used to train the network, the main disadvantage of which is premature convergence to the local optimum. In optimization using PSO, optimization variables in training a neural network include weights and biases related to the network. The work process is as follows: first, N location vectors X_i , where N is equal to the number of team members, are randomly generated. The group population is usually 4 to 5 times the number of optimization variables. The neural network is formed by the weights and biases obtained from the variables of these vectors and the error obtained from each execution is considered as the fitness rate of the variable vector of that network. This process is repeated until the final convergence is achieved. The aim of the final convergence is to reach the optimal location vector (values of optimal weights and biases), in such a way that the training error is minimized. Therefore, the objective function that should be minimized in this optimization is the amount of prediction error. If the number of layers is 3 and the number of neurons in the hidden layer is 7, so the number of weights is the number of weights between the input and the hidden layer + the number of weights between the hidden layer and the output (the number of inputs * the number of neurons in the hidden layer + the

number of neurons hidden layer = $10*7+7=77$) and the number of biases is the total number of neurons, which are 8. Therefore, the total number of decision variables is 85 and each group of population includes 5 vectors with this dimension. The purpose of standardizing the input data of the neural network model is to normalize them. The best situation for neural networks is when all inputs and outputs are between zero and one. To normalize the data, the mapminmax function was used in MATLAB software. The next step includes determining the model, specifying the architecture, the optimal number of iterations, determining the number of hidden and input layer neurons, the number of layers, and determining the appropriate stimulus function for the desired neural network model (in this research, the optimal number of iterations is 27,000, the number of layers is 3 and sigmoid tangent driving function is used). The next step is to train the network, which means determining the amount of weights and biases with a part of the data by optimization algorithms of particle swarm and genetics. Then the evaluation and testing of the network with the rest of the data is done by statistical parameters to evaluate the used algorithms and finally the output and simulation results are displayed by the model.

In this research, 85% of the data were considered for training and 15% of the data were considered for validating the model.

2.5. Model evaluation criteria: To determine the accuracy of the models the values of Root Mean Square Error ($RMSE$), Mean Absolute Error (MAE) and Determination Coefficient (R^2) was used:

$$RMSE = \sqrt{\frac{1}{n} \sum (y_{observed} - y_{predicted})^2} \quad 5$$

$$MAE = 100 * \frac{1}{n} \sum |y_{observed} - y_{predicted}| \quad 6$$

$$R^2 = 1 - \frac{\sum (y_{predicted} - y_{observed})}{\sum y_{predicted}^2 - \frac{y_{observed}}{n}} \quad 7$$

In the above equation, $y_{predicted}$, $y_{observed}$ and n are respectively the representatives of predicted values, observed values and the number of data. As the accuracy of the prediction model increases, $RMSE$ and MAE tend to zero and R^2 approaches one.

3. Results and discussion

In this study, quality data of Behbahan plain during the years 2014 to 2018 was used in order to simulate groundwater nitrate. Input information to the models for simulating nitrate (NO_3^-) on a monthly basis and including electrical conductivity (EC), calcium (Ca^{2+}), magnesium (Mg^{2+}), sulfate

(SO₄²⁻), bicarbonate (HCO₃⁻), chlorine (CL⁻), potassium (K⁺), hardness (TH) and pH (pH). Input

characteristics of quality parameters are shown in table (1).

Table 1 – Statistical Profile groundwater quality parameters Behbahan Plain

Water quality parameters	unit	minimum	maximum	average	Standard deviation
NO ₃ ⁻	mg/lit	1.28	112.4	24.95	17.41
EC	µs/cm	431	7457	2697.47	1594.08
Ca ²⁺	mg/lit	1.21	42.01	13.83	9.88
Mg ²⁺	mg/lit	0.5	28.51	6.69	5.2
SO ₄ ²⁻	mg/lit	0.79	47.86	26.88	12.96
HCO ₃ ⁻	mg/lit	1.04	5.88	3.47	0.9
CL ⁻	mg/lit	0.4	34.55	18.26	6.03
K ⁺	mg/lit	0.01	2.01	0.133	0.174
TH	mg/lit	51.5	285.03	938.3	409.7
pH	...	6.05	8.1	7.23	0.35

According to table 1, the highest and lowest values of nitrate concentration in groundwater are 112.4 and 1.28, respectively. While the maximum permissible and desirable nitrate in drinking water according to the national standard of Iran (No. 1053) is equal to 50 mg/lit [22]. The most important factor of nitrate entering this plain is the agricultural activities above the aquifer and the significant use of ammonia and potash nitrate fertilizers. In the flood irrigation method, about 20% of the total amount of irrigated water is returned to the groundwater, which is used in many areas of Behbahan plain. Also, the nitrogen compounds of urban and rural wastewaters are another factors of the high nitrate concentration in this plain. The high concentration of nitrate is due to the nitrification process [23]. Nitrogenous compounds enter the aquifer through urban and rural sewage and are converted into nitrates during the nitrification process. The high permeability of the sandy aquifer provides favorable conditions for the vertical transfer of oxygen to the deep parts of

the aquifer [24]. Due to the presence of sufficient oxygen, the nitrification reaction is possible, and most of the studied area is alluvial, and in alluvial aquifers with abundant oxygen, nitrification takes place and ammonia is converted into nitrate [16]. The mentioned reasons are consistent with the researches who conducted studies on nitrate concentration in drinking water of Tehran and Ardabil [25,26]. Due to the importance of knowing the state of nitrate concentration in the future, the artificial neural network model with two training algorithms for particle swarm optimization and genetics has been used to simulate nitrate concentration. All the calculations of this research were done in MATLAB, Excel and SPSS.

Analyzing the degree of correlation between the input variables and the work target variable is very valuable, because it provides useful information about the dependence of each of the input parameters on the target parameter. Table 2 shows the correlation between nitrate and other input parameters of the model.

Table 2 – The correlation matrix between nitrate and model input variables

Parameter	pH	TH	K ⁺	CL ⁻	HCO ₃ ⁻	SO ₄ ²⁻	Mg ²⁺	Ca ²⁺	EC
NO ₃ ⁻	-0.39	0.23	0.889**	0.157	0.51**	0.41	0.15	0.68**	0.72**

According to table 2, among all the input parameters of the model, the highest correlation with nitrate at the confidence level of 0.89 is related to K⁺ due to the use of potash fertilizers, which is based on the results of the research in [27] Nitrate of Zidon plain corresponds. In addition to K⁺, EC and Ca²⁺ have shown a high correlation at the confidence level of 0.99 compared to nitrate.

Different networks were formed with different number of neurons and repetitions. To avoid excessive learning and reduce errors, the optimal number of repetitions in the neural network model should be determined using evolutionary algorithms. To do this, the number of repetitions was changed from 5000 to 30000 and at the end of each step, its error was recorded. In the designed network, the error until 27,000 iterations showed a

downward trend and then an upward trend. Therefore, the optimal number of repetitions was chosen as 27,000. In the designed neural network model, the iteration that achieved the highest explanation coefficient in the training phase and the error below 5% was selected as the optimal iteration. Then, the number of different neurons for the hidden and input layers was considered to select the best model with the least error. The number of neurons in the hidden layer was changed from one to ten and the number of neurons in the input layer was changed from one to five, and in each step, the value of the RMSE and R² coefficients between the points simulated by the two algorithms and measured were calculated. The results of the repetitions of the neurons in the hidden layer are presented in table (3).

Table 3 – Results Repeat with different numbers of neurons in the hidden layer

Number of neurons in hidden layer	PSO		GA	
	RMSE	R ²	RMSE	R ²
1	0.816	0.941	0.767	0.847
2	0.828	0.942	0.789	0.861
3*	0.884	0.94	0.787	0.86
4	0.825	0.948	0.791	0.857
5	0.865	0.945	0.792	0.846
6	0.83	0.95	0.795	0.85
7*	0.811	0.957	0.794	0.869
8	0.815	0.956	0.794	0.874
9	0.822	0.938	0.798	0.879
10	0.848	0.939	0.802	0.883

According to table 3, the optimal number of neurons in the hidden layer for the optimization algorithms of PSO and GA are equal to 7 and 3, respectively. All the above steps were carried out to select the number of neurons of the input layer with the specified number of hidden layers, and the lowest error in the o PSO and GA was estimated with the number of input neurons 4.

RMSE, MAE and R² values were calculated between the points simulated by two algorithms

Table 4 – Statistics calculated during the training phase

Parameter	PSO			GA		
	RMSE	MAE	R ²	RMSE	MAE	R ²
NO ₃ ⁻	0.06	0.2	0.997	0.11	0.97	0.992

Table 5 – Statistics are calculated between the measured and simulated for calibration

Parameter	PSO			GA		
	RMSE	MAE	R ²	RMSE	MAE	R ²
NO ₃ ⁻	0.09	0.74	0.989	0.39	1.07	0.971

According to table (4) and (5), the highest accuracy in the simulation of underground water nitrate in Behbahan Plain is related to the particle swarm optimization algorithm, so that the RMSE and MAE values are the lowest and the R² index is the highest compared to the genetic algorithm in both the training and it is also in the calibration phase. PSO operates on the basis of search, and in this case, it can to some extent try to randomize the selection of points with a non-deterministic method. In this method, the particles try to select

(with optimal structure) and measured. Tables 4 and 5 show the statistics calculated in the training and calibration phase of the model, respectively.

and update the most optimal points in each iteration according to their positions and velocities, and for this reason, the simulation results are presented accurately. Among the other advantages of this algorithm compared to the genetic algorithm, we can mention easy implementation, low parameters of the algorithm and high convergence speed. Figure 3 shows the fitting of the curve between the measured and simulated points of the titrate concentration using two optimization algorithms of particle swarm and genetics.

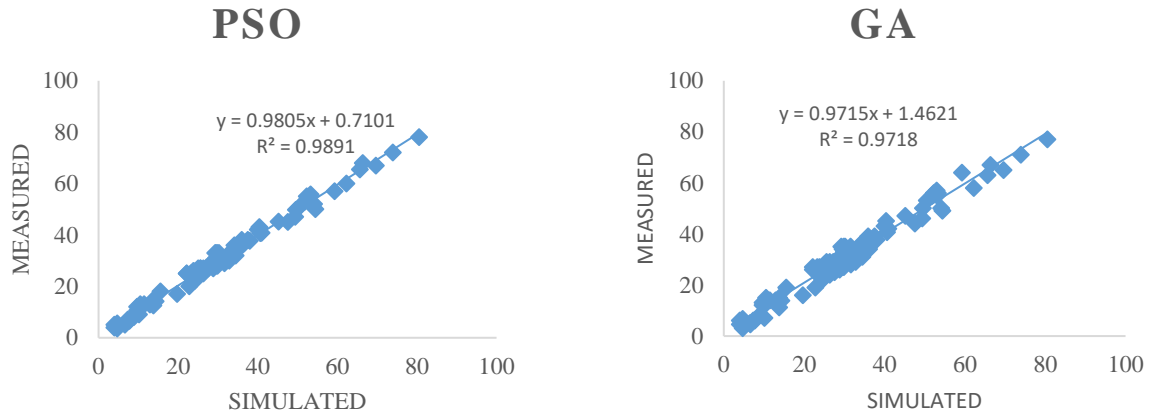


Figure 3- Scatter diagram measured and simulated nitrate concentrations using PSO and GA

In the graphs of the above figure, the R^2 value and the fitting equation between the points are specified. As mentioned, the high value of the R^2 coefficient indicates that the estimated values are close to the measured values. The linear fitting equation for each graph is defined as $y=ax+b$. The more the coefficient a tends towards one and the smaller the distance from one is, it indicates the better performance of the corresponding model,

which according to Figure 3, the value of the coefficient a in the particle swarm optimization algorithm is higher and as a result, it performs better than the genetic algorithm in this study. For the statistical comparison between the measured and simulated values of both algorithms, the statistical population mean comparison test was used using the t method at the error level of one percent, and the results are shown in Table 5.

Table 6 – The results of tests comparing the average

comparison	Measured and PSO+ANN			Measured and GA+ANN		
	MEAN DIFF	STD ERROR DIFF	P-value	MEAN DIFF	STD ERROR DIFF	P-value
NO_3^-	0.008	0.039	0.889 ^{n.s}	0.029	0.042	0.707 ^{n.s}

• n.s: there is no significant difference

Table 5 shows that the results of the GA have acceptable results for the simulation of nitrate concentration. So that there is no significant difference between the simulation values and the measured data of both algorithms at the error level of one percent.

Figures 4 and 5 show the comparison of the simulation of nitrate concentration with respect to time during the model testing period using the two used algorithms.

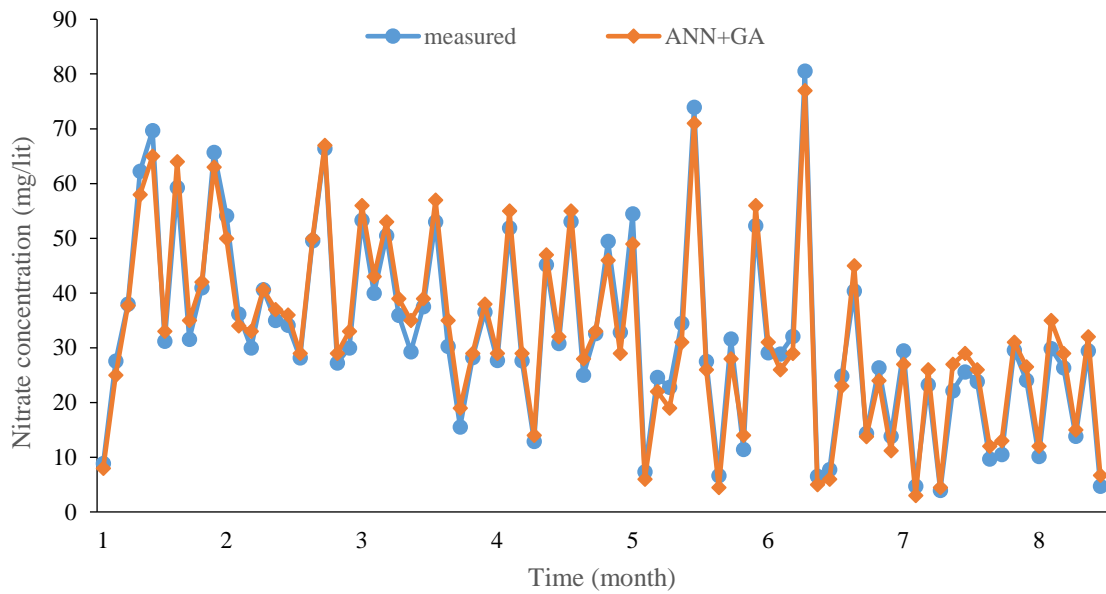


Figure 4- Compare simulate nitrate using a GA and measured data

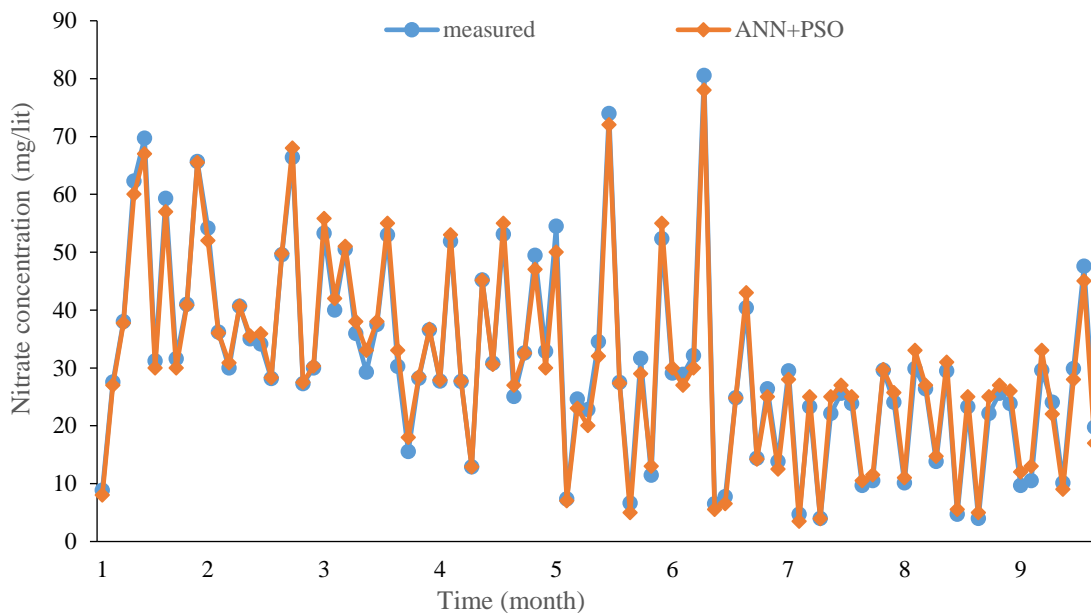


Figure 5- Compare simulate nitrate particles using PSO algorithm and measured data

4. conclusion

Artificial neural network is a suitable tool for adaptation, learning and classification of information. Many researchers have a great desire to use this tool, but they face the challenge of training neural networks. The combination of two ideas of collective intelligence and artificial neural network can be considered as an answer to this challenge. In this research, two optimization algorithms of particle swarm optimization and genetics have been used to predict nitrate concentration in Behbahan plain. The results showed that the accuracy of the particle swarm optimization algorithm is higher than the GA model. So that the values of RMSE, MAE and R²

statistics in the training phase for the PSO algorithm were equal to 0.06 (mg/lit), 0.2 (mg/lit) and 0.997, respectively. These statistics for the GA algorithm were determined as 0.11 (mg/lit), 0.97 (mg/lit), and 0.992, respectively. In the calibration stage, the RSME parameter for PSO and GA algorithms was calculated as 0.09 and 0.39 (mg/lit), respectively. MAE statistics for these two algorithms were estimated as 0.74 and 1.07 (mg/lit), respectively, and R² statistics were estimated as 0.989 and 0.971, respectively. The R² statistic in the simulation stage for PSO and GA algorithms was 0.989 and 0.971, respectively. Also, the results of the statistical test comparing the averages between the measured and simulated data

show that there is no significant difference between any of the values predicted by the used algorithms and the measured data; So these models can be used to determine nitrate concentration in groundwater sources.

References

- [1] Clark, I. and Fritz, P. Environmental Isotopes in Hydrogeology. 1st Edition. CRC Press. P. 300-342, 1997.
- [2] Khodaei, K., Mohammadzadeh, H., Naseri, H.R. and Shahsavari, A.A. Investigating nitrate pollution in the Dezful-Andimeshk plain and determining the source of pollution using N15 and O18 isotopes. Iranian Geology Quarterly, vol.22: p. 93-111., 2012(persian).
- [3] Keeney D., Nitrate in ground water: Agricultural contribution and control. Proceedings of the Agricultural Impacts on Ground Water Conference. p. 329-351, 1986.
- [4] Goulding, K. Nitrate leaching from arable and horticultural land. Soil Use and Management. Vol.16, no.1, p. 145–151, 2000.
- [5] U.S. EPA. Drinking water standards. U.S. EPA. 816-F-00- 995, p. 4, 1995.
- [6] Najah, A., Elshafie, A., Karim, OA. and Jaffar, O., Prediction of Johor river water quality parameters using artificial neural networks. European Journal of Scientific Research, vol. 28, no. 3, p. 422-435, 2009.
- [7] Asadollahfardi, A., Taklifi Gh. and Ghanbari A., Application of artificial neural network to predict TDS in Talkheh Rud River. Journal of Irrigation and Drainage Engineering, vol. 138, no. 4, p. 363–370, 2012.
- [8] Musavi-Jahromi, SH. and Golabi M., Application of artificial neural networks in the river water quality modeling: Karoon river, Journal of Applied Sciences, vol. 8, no. 12, p. 2324-2328, 2008.
- [9] Ramasamy, N., Krishnan, P., Bemard, J. and Ritter, W., Modeling Nitrate Concentration in Ground Water Using Regression and Neural Networks. FREC Research Reports, Department of Food and Resource Economics University of Delaware, 2003.
- [10] Yesilnacar, M.I., Sahinkaya E., Naz M. and Ozkaya B., Neural network prediction of nitrate in groundwater of Harran Plain, Turkey, Environ Geol, vol. 56, p.19–25, 2008.
- [11] Banejad, H., Kamali, Mahsa., Amirmoradi, Kimia. and Olyaie, F., Forecasting Some of the Qualitative Parameters of Rivers Using Wavelet Artificial Neural Network Hybrid (W-ANN) Model (Case of study: Jajroud River of Tehran and Gharaso River of Kermanshah), Journal Health & Environ, vol. 6, no. 3, p. 277-294, 2013(persian).
- [12] Mirzavand, M., Sadati Nrijad, M. and Akbari, M., Simulation Changes in groundwater quality with artificial neural network model (Case study: Kashan aquifer). Iranian Journal of Natural Resources, vol. 68. No.1, p. 159-171, 2015 (persian).
- [13] Shiri, J., and Kisi, O., Comparison of genetic programming with neurofuzzy systems for predicting short-term water table depth fluctuations, Comput. Geosci, vol. 37, no.10, p. 1692-1701, 2011.
- [14] Traore, S. and Guven, A., Regional-specific numerical models of evapotranspiration using gene-expression programming interface in Sahel, Water Resour. Manag, vol. 26 no.15, p. 4367-4380, 2012.
- [15] Moashrei, S.A., Tabatabaie, S.M., Razaghi, P., Sarani, N., and Eslami Mahdi Abadi, S.H., Estimating the groundwater nitrate by using artificial neural network and optimizing it by genetic algorithm., Proc. Environment and Civil Engineering, Kuala Lumpur, Malaysia, p. 85-92, 2012.
- [16] Stuart, M.E., Gooddy, D.C., Bloomfield, J.P. and Williams, A.T., A review of the impact of climate change on future nitrate concentrations in groundwater of the UK. Science of the Total Environment, vol. 409, no.15, p. 2859-2873, 2011.
- [17] Sayadi Shahraki, A., Naseri, A.A. and Soltani Mohammadi, A., Simulation of hydraulic head using Particle Swarm Optimization Algorithm and Genetic Algorithm. (Case study: Debal khazaie sugarcane plantation). Journal of Water Resources Engineering. Vol. 12 no.43, p.14-24, 2020 (In Persian).
- [18] Eberhart, R. and Shi Y., Comparing inertia weights and constriction factors in particle swarm, in: Proceedings of the Congress on Evolutionary Computation, La Jolla. p. 84–88, 2000.
- [19] Goldberg, D.E., Genetic Algorithms in Search Optimization and Machine Learning. Addison-Wesley. 1st Edition. p. 1-77, 1989.
- [20] Dayhoff, J. E., Neural Network Principles. Prentice-Hall International, U.S.A, 1999.
- [21] Khanna, T., Foundation of neural networks. Addison-Wesley Publishing Company, U.S.A, 1990.

- [22] Iranian National Standard No. 1053: 1376, physical and chemical characteristics of drinking water. (presian)
- [23] Saha, L.C. and Kumar, S., Comparative quality of potable waters at Bhagalpur, India *Acta Hydrochim Hydrobiol*, vol. 18, no.4, p. 459–467, 1990.
- [24] Hamilton, P. A. and Helsel, D.R., Effects of agriculture on groundwater quality in five regions of the United States, *Ground Water*, vol. 33, no. 6, p. 217–226, 1995.
- [25] Alighadri, M., Hazrati, S., Sakhaeizadeh, A. and Soleimanpour, M., Nitrate concentration measurement in Ardabil drinking water supply sources and distribution network, *Ardabil Health and Health Journal*, vol.2 no.2, p. 69-75, 2014 (presian).
- [26] Yazdanbakhsh, A.R., Mohammadi, H., Sheikhmohammadi, A., Bonyadinejad, R. and Ghanbari, Gh., Investigating the concentration of nitrite and nitrate in the drinking water of Tehran in the areas covered by Shahid Beheshti University of Medical Sciences, the 12th Iran Environmental Health Conference. Shahid Beheshti University of Medical Sciences, 2009. (presian)
- [27] Fazeli, M., Kalantari, N., Rahimi, M.H. and Khoubyari, A., Investigating the temporal and spatial distribution of Nitrate contamination of underground water sources in the Zidon Plain, *Water Resources Engineering Journal*, vol. 4, p. 45-51, 2011 (presian).



Paper Type (Research paper)

Transformer-based Meme-sensitive Cross-modal Sentiment Analysis Using Visual-Textual Data in Social Media

Zahra Pakdaman¹, Abbas Koochari^{1*} and Arash Sharifi¹

1. Department of Computer Engineering, Science and Research Branch, Islamic Azad University, Tehran, Iran.

Article Info

Article History:

Received: 2024/10/31

Revised: 2024/12/01

Accepted: 2024/12/08

DOI:

Keywords:

Visual Sentiment Analysis,

Textual Sentiment Analysis,

Vision Transformer, LDA,

SBERT Bi-encoder

*Corresponding Author's Email
Address: koochari@srbiau.ac.ir

Abstract

Analyzing the sentiment of the social media data plays a crucial role in understanding users' intentions, opinions, and behaviors. Given the extensive diversity of published content (i.e., image, text, audio, and video) leveraging this variety can significantly enhance the accuracy of sentiment analysis models. This study introduces a novel meme-sensitive cross-modal architecture designed to analyze users' emotions by integrating visual and textual data. The proposed approach distinguishes itself by its capability to identify memes within image datasets, an essential step in recognizing context-rich and sentiment-driven visual content. The research methodology involves detecting memes and separating them from standard images. From memes, embedded text is extracted and combined with user-generated captions, forming a unified textual input. Advanced feature extraction techniques are then applied: Vision Transformer (ViT) is employed for extracting visual features, while SBERT Bi-encoder is utilized to obtain meaningful textual embeddings. To address the challenges posed by high-dimensional data, Linear Discriminant Analysis (LDA) is used to reduce feature dimensionality while preserving critical classification information. A carefully designed neural network, consisting of two fully connected layers, processes the fused feature vector to predict sentiment classes. Experimental evaluation demonstrates the efficiency of the proposed method, achieving up to 90% accuracy on the MVSA-Single dataset and 80% accuracy on the MVSA-Multiple dataset. These results underscore the model's ability to outperform existing state-of-the-art approaches in cross-modal sentiment analysis. This study highlights the importance of integrating meme recognition and multi-modal feature extraction for improving sentiment analysis, paving the way for feature research in this domain.

1. Introduction

Sentiment analysis is a method to extract the user's real feelings from the data (text, image, video, and audio) that he/she publishes on the web platform. Considering that until the last few years, textual content was the most popular and most frequently published content on websites and social media, textual data attracted the attention of researchers. As a result, several methods were

presented in text sentiment analysis, and the methods proposed in this field have become very rich and fruitful.

In recent years, when image-oriented platforms such as Facebook, Instagram, etc., came to light, a large amount of image data was published by users. In this respect, the importance of analyzing the sentiment of this type of data was

also clarified for researchers. Therefore, with the appearance of convolutional neural networks (CNN), it becomes possible to provide more accurate methods of visual sentiment analysis. Yet, visual sentiment analysis faces certain difficulties. For instance, individuals often employ memes on platforms like Twitter to convey concepts. The utilization of this content has experienced a notable rise in recent years [2]. Many memes possess an ironic quality that cannot be accurately and fully captured by computer vision methods (Samples of the ironic meme are shown in Figure 1.). To address this issue, the article employs meme detection. Consequently, both the image and the embedded text undergo analysis, along with the user's caption. Moreover, sometimes the published text is not complete, expressive, and clear (i.e., the text only includes hashtags) to be able to understand the user's real sentiment. Therefore, to provide a more accurate analysis, cross-modal approaches – methods that use both text and image- are used to analyze the sentiment of social media content.

Social media images may include a wide range of topics (more than just a picture of a face or nature) which is a substantial challenge in visual sentiment analysis. As an efficient solution, the analysis of the published text along with the image will help to provide a more accurate analysis of the user's sentiment. This article presents a cross-modal approach that analyzes social media data with a wide range of topics. Also, as stated above, most of the bi-modal or cross-modal schemes -that use image and text for sentiment analysis- analyze different content in parallel; thus, the complexity of the architecture increases. However, the proposed method analyzes the content in an integrated way (combining the features from the beginning and obtaining the result of the sentiment analysis of those features in an integrated manner). Therefore, it will significantly reduce the dimension of the feature vector and the complexity of the architecture while accuracy will impressively increase.

In this article, to assess the sentiment of the image with the assistance of auxiliary data, the initial step involves determining whether the input image is a regular image or a meme. In memes, the embedded text is extracted and concatenated to the published text. Subsequently, for all input data, both image and text features are extracted. then, the dimension reduction method and the resultant vectors are concatenated. Ultimately, after passing the vector through the fully connected layers, the final sentiment label is

predicted. The main contributions of the paper are highlighted as follows:

For more precise analysis, a meme detector is employed to differentiate memes from other content. Hence, the embedded text in this content, which plays a crucial role in comprehending the genuine semantic and sentiment, will be extracted.

In this article, a transformer-based method is proposed. More accurate sentiment analysis is achieved thanks to the multi-head attention architecture used in the structure of transformers. This architecture makes it possible to capture both low-level features and high-level semantic information. Therefore, using transformers in both the image and text sides will help to solve the sentence ambiguity (like texts that only contain hashtags) and grasp the sense of irony or sarcasm of the image or text to predict the final sentiment.

Given to integrated architecture of the proposed method, the visual and textual features are combined from scratch, followed by a unified sentiment analysis operation on the concatenated feature vector.

In order to reduce the dimensionality of the extracted features, a dimension reducer (LDA) has been used. It improves the performance by better separation of sentiment categories.

In section 0 some sentiment analysis algorithms are reviewed. The proposed method is explained in section **Error! Reference source not found.** In section 0, the results of the proposed method are reported compared to state-of-the-art methods in this field. Finally, the conclusion of this paper will be presented in section 0.

2. Related works

Please be sure your sentences are complete and that there is continuity within your paragraphs. Check the numbering of your graphics and make sure that all proper references are included.

Various categories are presented for different types of content sentiment analysis schemes. One of the introduced taxonomies is based on the number of contents used in the analysis process. Based on this classification criteria, existing methods are divided into unimodal, bi-modal, and multi-modal. In unimodal methods, only one content is analyzed for sentiment analysis. In the bi-modal (cross-modal), two types of content are used, and in multi-modal ones, more than two contents are used to estimate the final sentiment. Since a cross-modal approach is presented in this article, cross and multi-modal methods will be reviewed in this section.



Figure 1 Examples of ironic memes. (a) and (b) embedded text changes the sentiment of images. (c) the image has both positive and negative content.

2.1. Cross-modal approaches

The published data in social media includes a wide range of topics, including person’s faces, landscapes, memes, abstract images, etc., which turns “analyzing their sentiments” into a challenge. Moreover, the presence of sarcasm and irony in viral images and memes on social media can create another challenge in analyzing accurate sentiment. One of the ways to improve the results of visual sentiment analysis is to use the auxiliary data (often text) published along with images in social media. The use of auxiliary data can provide a more comprehensive result regarding the genuine sentiment of the published data. These methods are called bi-modal (cross-modal).

Zhang et al. [2] presented a multimodal method in which CNN and GloVe are used to extract image and text features, respectively. The extracted features are given to IDLSTM as an image-text pair to finally obtain the result of sentiment analysis. Lu et al. [4] used BERT to process text and image in parallel. Since BERT was proposed for text analysis, in their article, by applying changes to this network, it is used for text and image analysis simultaneously. They used two BERTs in parallel. BERTs process regions of image and pieces of text. The difference between the BERT structure in this paper and the original BERT is that in their paper, a co-attention transformer layer is used to cross-feed image regions and text pieces to the Visual and linguistic transformer. Xu et al. [5] believe that the relationship between the contents of image and text is a bidirectional one. For this reason, they used text and image together. In this way, they can report more accurate results for sentiment analysis. The authors proposed a model called Bi-directional Multi-level Attention. In the architecture of this method, two attention

networks have been used in parallel form. The complex relationship between the image and its corresponding text – which includes a description of that image- will be analyzed. At last, the final label will be predicted based on the obtained information from the image and text. In the proposed method by Huang et al. [6] sentiment analysis will be done through three parallel axes. In their approach, the sentiment labels of the text, image, and text-image are predicted separately. Eventually, the final label will be obtained in a late fusion way. Li et al. [7] extracted features of the image and translated the image into text using LSTM. Then, using the SentiBank dataset, a number of Adjective-Noun Pairs (ANP) will be generated for each image and the generated caption will be combined with ANPs. Finally, the sentiment of the obtained text will be predicted by a bi-directional LSTM. The proposed method by Serra et al. [8] used a teacher-student architecture which will bring about the textual data sentiment using a pre-trained BERT model. In the student section, ViT is used to obtain the sentiment label of the image data. Eventually, by setting a threshold for the loss function, the final sentiment label will be obtained. Zhu et al. [9] have focused on extracting relevant regions of text and image. First, a pre-trained Faster R-CNN is used to recognize image regions and their corresponding representation. In order to embed each word in a vector and summarize the text information in the sentence, a pre-trained BERT and a bi-directional GRU are used, respectively. Also, to extract information from the image and relevant features in both text and image, pre-trained ResNet18 and cross-modal attention are employed, respectively. Finally, the final sentiment is predicted by concatenating the features and passing through Softmax. Yadav et al. [10] extract image features by two types of attentions (channel and spatial

attentions) and text features by LSTM and GRU. Then, the final label is estimated using the attention mechanism. In [11], the visual and textual features are extracted by pre-trained BERT. Then, the sentiment of each data is obtained by two contrastive learning structures (label and data-based contrastive learning). The proposed scheme by Peng et al. [12] consists of three modules. The first module is the feature extraction module, which extracts text and image features. The second is the feature attention module, which is responsible for learning text and image attention features. The third is the cross-modal hierarchical fusion module, which combines the features of within and between modals. In [13], Xu and Mao used two pre-trained Object-VGG and Scene-VGG networks to extract object and scene features of the image, respectively. To extract text features, they used LSTM. At last, the sentiment label is predicted by passing the extracted features through fusion and Softmax layers.

2.2. Multi-modal approaches

Most multimodal sentiment analysis methods deal with datasets of speech videos of different people. Therefore, unlike the online published content, which covers a wide range of topics, they have a limited scope. In these methods, the person's speech text, audio, and video have been used to provide the final sentiment label. In such datasets, real emotion can be easily obtained by analyzing the tone of voice, laughter, crying, facial expression, etc.

Wen et al. [14] presented a method in which text, audio, and video features are first extracted by GloVe, COVAREP, and facet respectively. In order to modify the features, attention structure is used. After merging and modifying the extracted features two by two, the final label will be predicted using an attention transformer and a fully connected layer. This scheme has only been tested to analyze the sentiment of people's speech. Yang et al. [15] proposed a method that tried to analyze the sentiment of the text more precisely by using audio and video. In their method, using GRU, audio and video are translated into text and analyzed together with the existing text. In fact, like using Seq2Seq methods – which consist of an encoder and a decoder part- one language can be translated to another; these methods can be used to translate audio or video into text. This approach is also limited to the sentiment analysis of people's speech. The authors used an encoder-decoder architecture to design their proposed model. The proposed method by Rahman et

al. [16] is a BERT-based [17] method, which processes text, audio, and video altogether through a mechanism called Multimodal Adaptation Gate (MAG) among BERT layers. This approach used two BERT architectures in a row to analyze people's feelings during the speech. Hazarika et al.' method [18] used text and audio as auxiliary content for more detailed speech video analysis. To increase the accuracy of the analysis, two categories of features are extracted; common and specific ones. By defining and extracting these two categories of features, the model will be more comprehensive because both specific and common features of different contents are considered. They used a pre-trained BERT model to extract the text feature and an sLSTM to extract the audio and video features. Mittal et al. [19] presented another multimodal method for speech sentiment analysis, the highlight of which is reducing the impact of noisy data and being able to make decision in the presence of noise. In this method, audio, video, and text features are extracted. If there is noise in the extracted features, the features are changed to noise-free data. Also, effective features are separated from ineffective ones so that effective ones can be used. The audio, video, and text features are individually fed into two-layer networks. Then the output of these networks is given to an attention layer. Finally, after several fully connected layers, the result of sentiment analysis will be predicted.

3. Proposed method

As Figure 2 shows, to enhance the accuracy of predicting the sentiment class of the image, a cross-modal approach is introduced, incorporating the published text alongside the image. initially, a binary classifier is applied to identify memes. Memes need to be preprocessed before entering the proposed architecture. Then, visual and textual features are extracted using the transformers. To reduce dimensionality and enhance classification, the dimensions of the features are reduced by Linear Discriminant Analysis (LDA). Ultimately, the reduced feature vectors are concatenated, and by passing through two fully connected layers, their sentiment label will be estimated. In the following, the steps of the proposed method are explained in detail.

3.1. Meme detection

Since both the image and embedded text play crucial roles in the meme's sentiment label prediction, and a substantial portion of published images on social media constitute memes, this

article includes a preprocessed section dedicated

to meme detection and extracting embedded text.

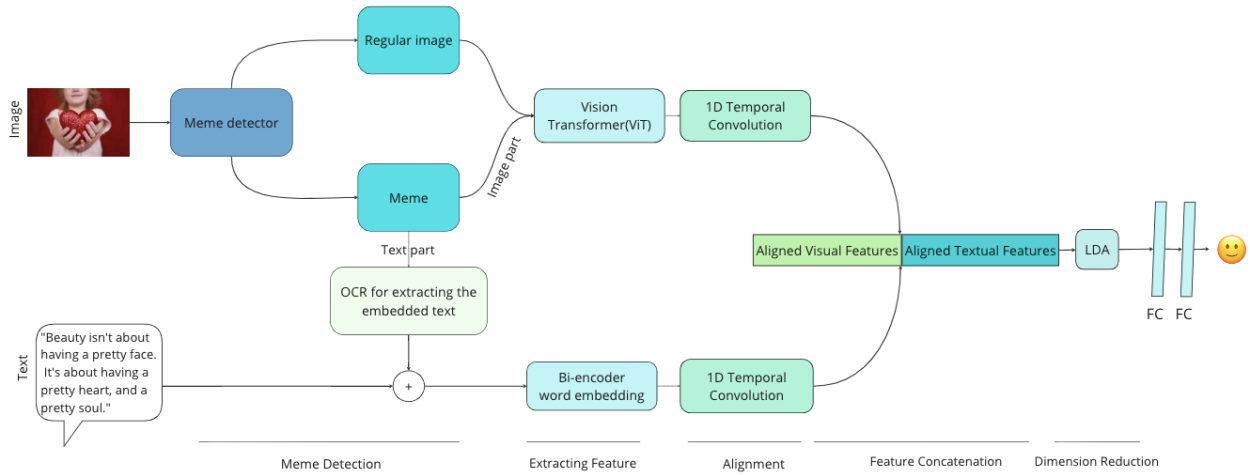


Figure 2 The proposed method flowchart, the final label is predicted after extracting features, reducing dimensionality, and passing through two fully connected layers.

To differentiate between a meme and a regular image, a binary classifier is employed. This classifier utilizes a pre-trained Xception as a base model, followed by three fully connected layers with 128, 64, and 32 hidden units to classify an input image into either a meme or regular class. The network is trained using a meme dataset introduced in section 0.

As depicted in Figure 2, if an image is identified as a meme, the embedded text is retrieved using Optical Character Recognition (OCR). The retrieved text part is then concatenated to the user’s caption, while the image part is utilized to extract features. The recognition and extraction of embedded text in a meme are performed using EasyOCR module, an open-source tool known for its speed and efficiency in word recognition within images, along with its adaptability in handling various typeface and text layouts.

3.2. Feature extraction and alignment

In the second step, the text and image features are calculated using the transformer. Visual Transformer (ViT) and SBERT Bi-encoder are used respectively to extract image and text features. Then, in order to align these features, 1D temporal convolution will be applied to the text and image features separately. How to extract features using ViT and Bi-encoder are explained below.

3.2.1 Image feature extraction by ViT

The main idea of the transformer is to use the attention architecture to process the data. The use of transformers has found a special place in natural language processing problems, and they

are used in almost all proposed methods. But Dosovitsky et al. [20] showed that this level of dependence on CNNs is not necessary in image processing applications and decided to use transformer in image processing areas.

The vision transformer learns a hierarchical representation of the image by stacking many layers of patch embeddings, multi-head self-attention, and feedforward networks. This capability makes it possible to capture both low-level features and high-level semantic information [21]. Using this method, tasks such as image classification can be performed well with only one transformer applied directly to a sequence of image patches. As Figure 3 shows, the overall architecture of this method is summarized in the following steps:

Step 1: The input image (x with resolution (H, W)) is divided into non-overlapping patches (x_p with resolution (P, P)). The number of patches is $N = HW/P^2$.

Step 2: The patches are flattened.

Step 3: From flattened patches, linear projection with lower dimensions is produced.

Step 4: Position embedding is added to each projected patch. Like $[class]$ token of BERT, a learnable embedding is prepended to the sequence of patches ($z_0^0 = x_{class}$).

Step 5: The generated sequence is given as input to a Transformer encoder.

In the encoder, after performing normalization (Layer Norm (LN)) on the embedded patches (flattened patches with position embedding), the output is given to the multi-head self-attention unit (MSA). In this unit, each patch produces three vectors. The first vector is Query, which is based on the patch’s meaning. Each patch creates this vector to announce its requirements to other

patches or to the model. The next vector is Key. Through this vector, each patch tells other patches what information it brings. The Value vector is the headline of a document, and based on this headline, it is decided how much information can be obtained from a document. After obtaining the Key, Value, and Query vectors, embedded patches are added to the result of MSA. This operation will be done in order to prevent the creation gradient vanishing problem and re-adding position embedding information after applying attention. Then, normalization will be done again, and the result will be passed through a Multi-Layer Perceptron (MLP). Finally, the MLP result will be added to the previous addition result. Passing through the Transformer encoder layer will be done L times. The operation of the Transformer encoder will be formulated as follows:

$$z_0 = [x_{class}; x_p^1 E; x_p^2 E; \dots; x_p^N E] + E_{pos} \quad (1)$$

$$z'_l = MSA(LN(z_{l-1})) + z_{l-1}, \quad l = 1, \dots, L \quad (2)$$

$$z_l = MLP(LN(z'_l)) + z'_l, \quad l = 1, \dots, L \quad (3)$$

$$y = LN(z_L^0) \quad (4)$$

Where z_L^0 is the output of the Transformer encoder and y is the image representation.

Step 6: In this step, the features can be extracted or the output class can be estimated through MLP.

3.2.2. Sentence embedding by Bi-encoder

For generating sentence embedding of the published text with the image, a Sentence BERT/roBERTa (SBERT) architecture, which is called Bi-encoder, is used. This sentence embedding method, as a transformer-based method, is much more efficient and faster than other similar ones [22]. Sentence embedding has several advantages over other common NLP tasks [23]:

The firstly mentioned (sentence embedding) is more comprehensive compared to other common NLP tasks, thanks to examining the whole sentences rather than word level feature extraction (common in other NLP approaches).

Since sentence embedding captures the meaning and context of a sentence in a numerical form, it can easily be used as an input to models and helps to understand the content of the text. Thus, performance will be improved.

Sentence embedding represents a sentence or a phrase as a single vector with a fixed length. So, the dimensionality of data will be reduced and it will be easier to work with.

The sentence embedding is resistant to the presence of noise and variability (such as spelling errors and variation in language usage). This property makes it suitable for some tasks like text classification and sentiment analysis, where the input sentences may not always be clear and uniform.

As Figure 4 shows, first, the input sentence (s_i) is tokenized: $T = [T_1, \dots, T_n]$. Tokens are fed to a BERT language model to obtain fixed-size word embeddings. As stated in the previous section, the attention mechanism has three qualities: Value (V), Key (K), and Query (Q). In SBERT architecture, the values of V , K , and Q for each token are calculated as follows:

$$v = W_v^T T_i \quad v \in \mathbb{R}^p \quad V = [v, \dots, v_n]_{p \times n} \quad (5)$$

$$k = W_k^T T_i \quad k \in \mathbb{R}^p \quad K = [k_1, \dots, k_n]_{p \times n} \quad (6)$$

$$q = W_q^T T_i \quad q \in \mathbb{R}^p \quad Q = [q_1, \dots, q_n]_{p \times n} \quad (7)$$

Where, W_v^T , W_k^T , and W_q^T are transformation matrices.

After calculating the values of V , K , and Q , MSA is applied to them. Passing through the normalization and feedforward layers, the transformer result will be generated. Finally, in order to reduce the dimension, it will be passed through a pooling layer.

$$Multihead(Q, K, V) \quad (8)$$

$$= W^0 \text{concat}(head_1, \dots, head_h)$$

Where, W^0 is a linear transformation. $head_i$ s are the output of different heads. H is the number of repeating blocks.

3.2.3. Feature alignment

After extracting the image and text features, to ensure that each element of the extracted feature has sufficient awareness of its neighboring elements, they are passed through a 1D temporal convolution layer, as follows [24]:

$$\hat{X}_{\{Im, TxT\}} = Conv1D(X_{\{Im, TxT\}}) \quad (9)$$

Where, $\{Im, TxT\}$ are visual and textual content, respectively. X_{Im} is the extracted feature of the image. X_{TxT} is the extracted feature of the text.

Therefore, it is expected that the convolved result contains the local information of the feature. In addition, the extracted features that have different dimensions will be projected to the same dimension by 1D temporal convolution.

3.3. Dimension reduction by LDA

In this stage of the proposed method, after extracting the features and performing the alignment, the resulting vectors are concatenated.

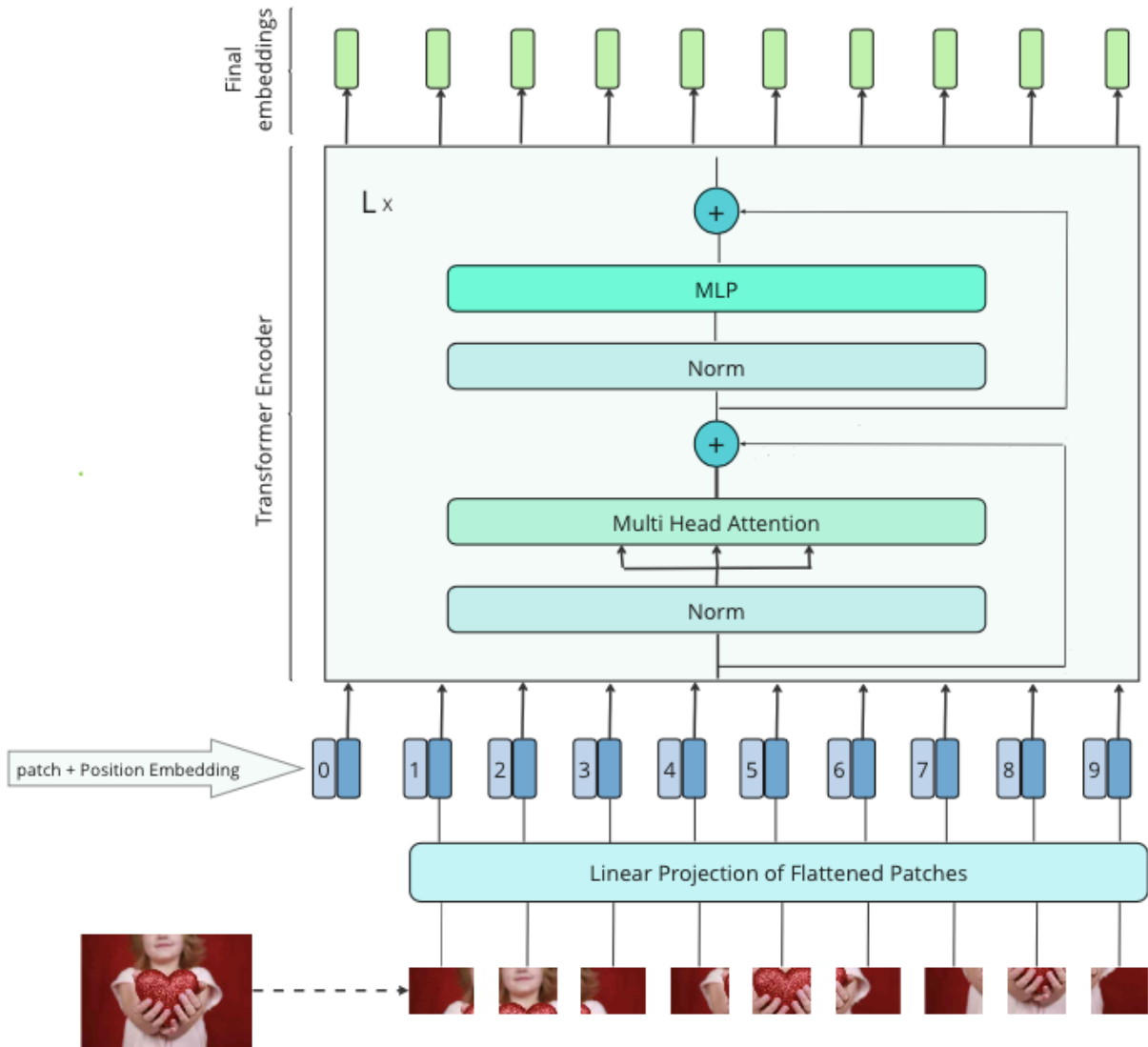


Figure 3 Visual feature extraction; Features of the input image are extracted by ViT. The input image is divided into non-overlapping patches. The patches are flattened. From flattened patches, linear projection with lower dimensions is produced. Position embedd

In order to improve the performance by better separation of sentiment categories, LDA is used to reduce the dimensionality of the concatenated vector.

The presence of high dimensions in the feature space can mean that the points in the space often represent very small and non-representative samples. This problem, which is known as the curse of dimensionality, can greatly affect the performance of the proposed method. This problem can often be solved by reducing the dimensions of the input features.

LDA is a predictive modeling algorithm used for multi-class classification. This algorithm can be used as a dimension reduction approach. This transformation projects the features from the higher dimensional space to the lower one while

maintaining the most important structure or relations between the observed variables.

LDA first calculates the inter-class variance to reduce the dimensions of the features. For this purpose, it calculates the distance between the averages of different classes in order to obtain the amount of separation between classes.

$$S_b = \sum_{i=1}^c N_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^T \tag{10}$$

Where c is the number of classes, N_i is the number of samples of class i , \bar{x}_i is the average of i -th class samples and \bar{x} is the average of all data points.

Then the intra-class variance is calculated. To do this, the distance between the average and sample of each class will be computed.

$$S_w = \sum_{k=1}^c S_k, S_k = \sum_{x(n) \in C_i} (x(n) - \bar{x}_i)(x(n) - \bar{x}_i)^T \quad (11)$$

Finally, LDA creates a space with lower dimensions that maximize the inter-class variance and minimize the intra-class variance. In this way, P is considered as the lower dimensional space.

$$P_{lda} = \arg \max_P \frac{|P^T S_b P|}{|P^T S_w P|} \quad (12)$$

According to what was discussed about the method and benefits of dimension reduction by LDA, in the proposed method, this approach is used to reduce the feature vector resulting the concatenated of image and text features.

3.4. Label estimation

After extracting the visual and textual features, aligning and concatenating the features, and reducing the dimensions of the result vector, the final emotion label should be estimated at this stage. Therefore, two fully connected layers have been used. The reduced vector is fed into these two layers, and the final label of the sentiment class is estimated.

4. Experimental results

In order to estimate the sentiment of social media content, a meme-sensitive cross-modal method is proposed in this article. According to this approach, visual and textual content are used to determine the final sentiment class label.

To implement this method, first, image and text features are extracted using the ViT model and SBERT Bi-encoder method, respectively. Then, in order to align the extracted feature vectors, 1D temporal Convolution was used. Aligned vectors are concatenated, and the resulting dimensionality is reduced by LDA transformation. Then, they have been passed through two fully connected layers with ReLU activation function, which have 64 and 32 hidden units, respectively. Eventually, in the last layer, the sentiment label is estimated by the Softmax activation. The implementation is run on an NVIDIA GeForce RTX 2060. Also, the proposed method will be evaluated through Precision, Recall, F1, and Accuracy measures.

4.1. Datasets

In examining the results of the proposed method and comparing it with the performance of other state-of-the-arts, the following datasets have been used:

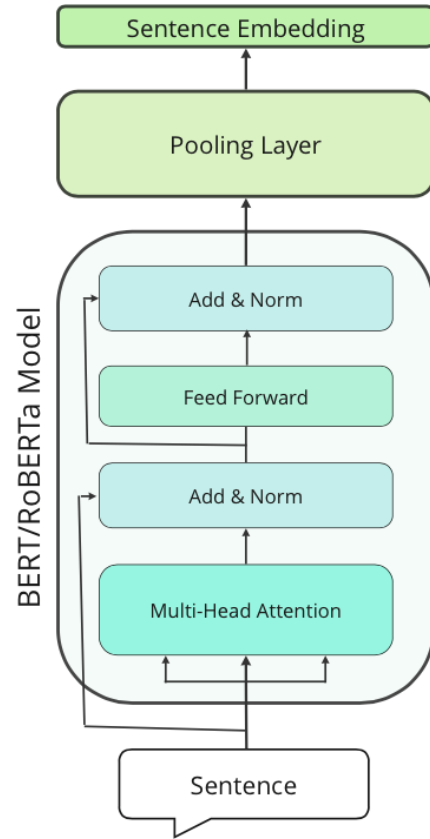


Figure 4 Sentence embedding

MVSA-Single and **MVSA-Multiple** [25]: In these datasets, text and image content of Twitter have been collected. The data are annotated into three categories: positive, negative, and neutral sentiment. The MVSA-Single containing 5,129 text-image pairs, is labeled by one annotator while the MVSA-Multiple with 19,600 text-image pairs is labeled by three annotators. In order to make the comparison fair, the content of both datasets was processed through the method proposed in [13]. According to this method, pairs that have different sentiment labels for text and image will be removed. Also, if one has a positive (or negative) label and the other has a neutral label, the positive (or negative) sentiment will be considered as the final label. Therefore, the number of pairs will be reduced to 4,511 in MVSA-Single and 17,024 in MVSA-Multiple.

Meme dataset: To collect a dataset for meme detection, the Hateful Memes dataset [26], consisting of 10,000 image memes, is utilized. This multimodal dataset is designed for detecting hateful memes, but the study does not consider the nature of these memes (whether they are hateful or not). All images from this dataset are treated as memes. Additionally, 10,000 random images from the COCO dataset are collected for regular image samples.

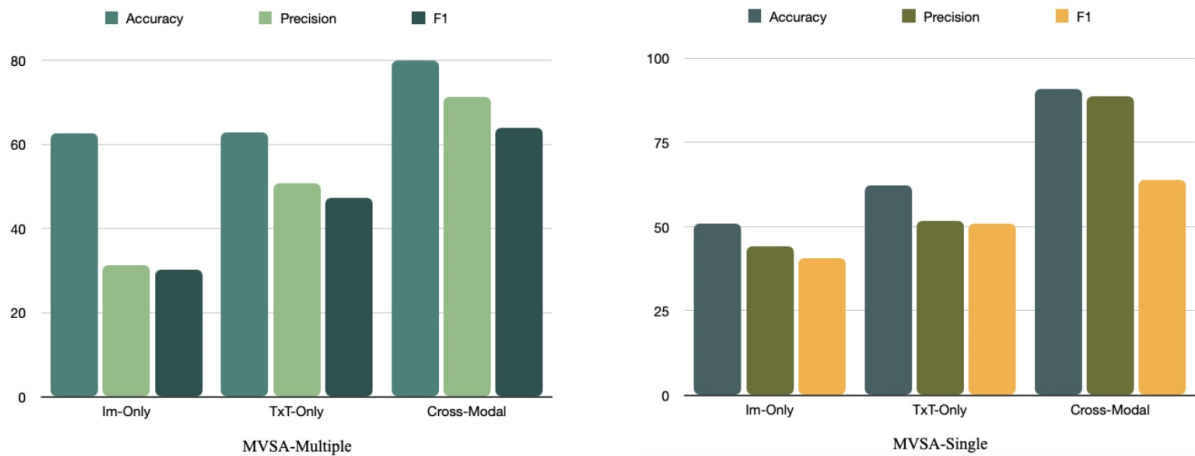


Figure 5 Results of comparison Accuracy between the Visual and Textual sentiment Analysis and the proposed method on MVSA datasets

4.2. The effect of using both image and text (cross-modal)

In this article, a cross-modal approach is proposed that uses visual and textual data in order to achieve a more accurate sentiment label. In this section, an experiment has been designed to prove that the use of text or image data alone has a significant effect in estimating the final result more accurately.

Im-Only: In this method, image features are first extracted using ViT. The dimensions of the features are reduced by LDA transformation and the final label is predicted after passing through two fully connected layers.

TxT-Only: First, sentence embeddings of input texts are calculated by SBERT RoBERTa model. After performing dimension reduction, the final text class is estimated using two fully connected layers.

As Figure 5 shows, on MVSA datasets, cross-modal has helped to more accurately estimate the sentiment of content in social media. This approach has a better performance than when only one type of content (text or image) is used.

4.3. The effect of using different feature extractors

This section will examine the effect of different feature extractors on textual and visual content.

4.3.1. Performance of different Visual feature extraction approaches

In the proposed method, ViT has been used to extract image features. Due to the fact that this approach is based on the transformer and its essence of attention architecture, it will extract better features and ultimately lead to a more accurate sentiment label prediction.

In this section, in order to compare the performance of different image feature extraction

schemes, a number of well-known feature extractors of mid and high-level features, will be examined. In order to extract mid-level features, SIFT, ORB, and LBP, and to extract high-level features, pre-trained VGG16 and ResNet50 have been used. As mentioned in literatures [27], features extracted from deep networks are more efficient than low and mid-level features. In Table 1 the comparison results of different feature extractors on MVSA datasets have been compared. As expected, the results of ViT feature extraction are much better than other methods, even the feature extraction by pre-trained CNNs. Since in this method, the attention mechanism is used, the extracted features have higher efficiency. Therefore, they will have a significant impact on the final estimation.

4.3.2. Performance of different textual feature extraction approaches

As stated in 0, the SBERT Bi-encoder method has been used to obtain sentence embedding. This feature extraction method based on Sentence BERT/RoBERTa has a much better performance compared to traditional word embedding methods. As mentioned in [22], BERT or RoBERTa model can be used in this approach. According to the results [22], they do not have significant differences. However, in this section, in order to prove the better performance of these two models (Sentence BERT/RoBERTa) than the other pre-trained models, the applying results are compared with other embedding models;

MiniLM model [28]: It is a sentence transformer model that maps sentences or paragraphs to the vector space. This model, which is pre-trained on more than one billion training data, provides decent quality while being fast.

Universal Sentence Encoder (USE) [29]: This model is also a transformer-based model that

Table 1 Results of the different visual feature extractors on MVSA datasets (In all cases, SBERT-RoBERTa method is used for textual feature extraction and LDA is used for dimensionality reduction.)

Methods	MVSA-Multiple				MVSA-Single			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
SIFT	0.6386	0.5098	0.4787	0.4707	0.7284	0.7045	0.6789	0.6731
LBP	0.6351	0.5139	0.4815	0.4724	0.7284	0.6561	0.6108	0.6112
ORB	0.6433	0.4558	0.4201	0.3960	0.6220	0.6257	0.6151	0.5757
CNN-VGG16	0.5628	0.4357	0.4404	0.4132	0.6685	0.6908	0.6565	0.6504
CNN-ResNet50	0.5622	0.4295	0.4314	0.4061	0.6763	0.7216	0.6744	0.6755
ViT	0.8043	0.7171	0.6484	0.8060	0.9089	0.8860	0.8527	0.9078

encodes sentences into the vector space using a transformer encoder and deep averaging network.

Doc2Vec [30]: this method is an unsupervised one for generating a fixed-length feature vector for a piece of text (like a sentence, paragraph, or document). This method is introduced as an extension of Word2Vec. In Doc2Vec, a three-layer deep neural network is used. This scheme is not a unified method like Bag Of Word, but has two different variations of Skip-Gram and CBOW.

The results of applying different pre-trained text feature extraction approaches on datasets are shown in Table 2. As can be seen, the transformer-based methods provide better results than the Doc2Vec. Among the transformer-based schemes, the SBERT feature extractor and the RoBERTa model gives the best performance. The RoBERTa is superior than BERT model due to being trained in larger datasets, using a more effective training procedure, and using a dynamic masking pattern instead of the static one used in BERT.

4.4. The effect of using different dimension reduction methods

In the proposed method, LDA transformation is used in order to reduce the dimension of the vector resulting from concatenation of image and text features. In this section, the efficiency of using or not using the dimensionality reduction method (LDA) is presented. As can be seen in Table 3, when LDA is used to reduce the dimensionality of feature vectors, the values of measures on MVSA datasets have grown significantly compared to not using dimension reduction.

PCA Transformation: It is an unsupervised linear transformation technique that identifies patterns in data based on correlations between features. In fact, PCA seeks to find the directions of maximum variance in high-dimensional data and projects it into a new subspace that has dimensions less than or equal to the original space.

SVD Transformation: This transformation is a popular method for dimensionality reduction. This method is a data-driven one because patterns are extracted from data without any expert knowledge or intuition. Computing SVD consists of finding the eigenvalues and eigenvectors. The SVD is a steady numerical transformation. It creates a representation of data based on dominant correlation within data.

In Figure 6, the ROC plots of different dimensionality reduction approaches on MVSA-Single are illustrated. As can be seen, the results of LDA transformation have been much better compared to the other two transformations.

4.5. Performance comparison

In this section, the proposed method has been compared with various schemes that have been proposed in the field of sentiment analysis in social media.

The comparison of the results of different methods on MVSA (Single and Multiple) is shown in Table 4. As it is clear, the proposed method has achieved significant results compared to other methods. In Xu's Method [13], text and image features are extracted by LSTM and pre-trained CNNs, respectively. Eventually, the sentiment label will be obtained through Softmax. Although Zhu et al. [9] also, have used CNN and GRU to extract features, they used a Faster R-CNN network to extract the corresponding features in both text and image. In the proposed method by Yadav et al. [10], Similar to previous methods, CNN and RNN were used for feature extraction. However, it used several attention structures to estimate the sentiment class with higher accuracy. The difference between Li's method [11] and the other aforementioned methods is the use of two contrastive learning structures to estimate the sentiment class. The three-module method of Peng et al. [12], in addition to extracting features by usual methods, uses attention structure and cross-modal hierarchical fusion to extract corresponding and specific features.

Table 2 Results of the different textual feature extractors on MVSA datasets (In all cases, ViT approach is used for visual feature extraction and LDA is used for dimensionality reduction.)

Methods	MVSA-Multiple				MVSA-Single			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
Doc2Vec	0.7462	0.5130	0.4398	0.6958	0.7537	0.7053	0.6832	0.7552
USE	0.7743	0.7289	0.7691	0.7294	0.8571	0.8515	0.8487	0.8562
MiniLM	0.7704	0.6408	0.5667	0.7283	0.8350	0.7848	0.7427	0.8357
SBERT-BERT	0.7926	0.6851	0.5990	0.7470	0.8892	0.8758	0.8538	0.8885
SBERT-RoBERTa	0.8043	0.7171	0.6484	0.8060	0.9089	0.8860	0.8527	0.9078

Table 3 Results of using and not using LDA on MVSA datasets (In all cases, ViT and SBERT-RoBERTa are used for visual and textual feature extraction, respectively.)

Datasets	Without LDA				With LDA			
	Accuracy	Precision	Recall	F1	Accuracy	Precision	Recall	F1
MVSA-Multiple	0.6204	0.4625	0.4378	0.6186	0.8043	0.7171	0.6484	0.8060
MVSA-Single	0.6053	0.5657	0.5174	0.6090	0.9089	0.8860	0.8527	0.9078

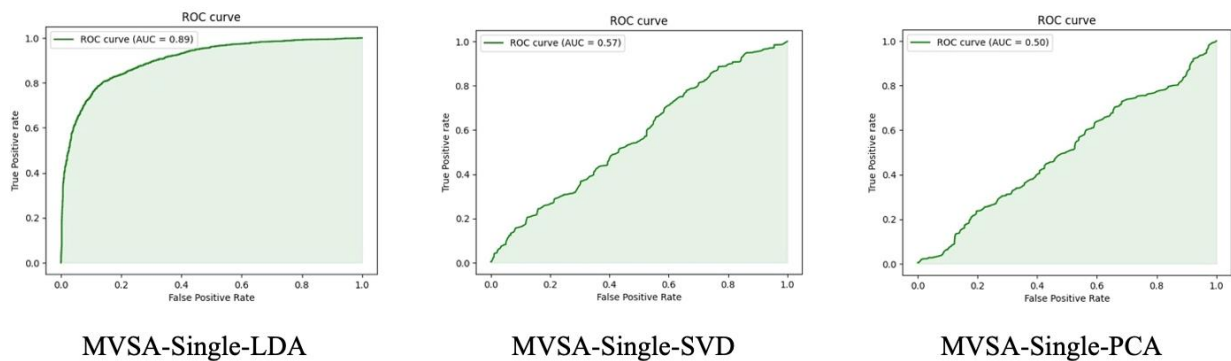


Figure 6 ROC plot - AUC of the different dimension reduction methods on MVSA-Single dataset

Table 4 Comparison between the results of the proposed method and other state-of-the-arts MVSA (Single-Multiple) datasets

Method	MVSA-Single		MVSA-Multiple	
	Accuracy	F1	Accuracy	F1
Xu's Method [13]	0.6984	0.6963	0.6886	0.6811
Zhu's Method [9]	0.7519	0.7497	0.7352	0.7349
Yadav's Method [10]	0.7959	0.7947	0.7526	0.7789
Li's Method [11]	0.7533	0.7346	0.7200	0.6983
Peng's Method [12]	0.7361	0.7503	0.7045	0.7477
Proposed Method	0.9089	0.9078	0.8043	0.8060

5. Conclusion and future works

Since the sentiment analysis of published content in social media is an abstractive issue and its content cover a wide range of topics (full of ironic, metaphorical, and abstractive concepts), using a single content may not lead us to the exact sentiment of the content. In this respect, the use of auxiliary data helps to predict a more accurate result. Therefore, in this article, a cross-modal approach is introduced to predict the sentiment class of published content. Rather than relying on

a single content type, both the image and the text published by the user are employed to determine the final label. Additionally, recognizing the growing prevalence of memes on social media, this method includes meme detection and preprocessing to ensure a more accurate analysis of the final sentiment. Since transformers have made significant progress in both visual and textual data analysis in recent years, this scheme has been used to extract features from pre-trained models based on transformers. Also, in the proposed method, LDA transformation has been

used to reduce the dimensionality of the extracted features and enhance classification. Finally, as reported in the experimental section, the performance of this method is much better compared to other state-of-the-arts. To improve the proposed method in the future, the corresponding relationships between visual and textual content can be used. In this case, it is possible to extract features that can express intra and extra-content relations in the best way. Eventually, the performance of the method can be increased, especially, for the analysis of content with ironic and metaphorical concepts.

References

- [1] Jain R, Rai RS, Jain S, Ahluwalia R, Gupta J (2023) Real time sentiment analysis of natural language using multimedia input. *Multimed Tools Appl.* <https://doi.org/10.1007/s11042-023-15213-3>
- [2] Ramamoorthy S, Gunti N, Mishra S, Suryavardan S, Reganti A, Patwa P, Das A, Chakraborty T, Sheth A, Ekbal A, Ahuja C (2022) Memotion2: Dataset on sentiment and emotion analysis Memes. In *Proceedings of De-Factify: Workshop on Multimodal Fact Checking and Hate Speech Detection*, CEUR.
- [3] Zhang K, Zhu Y, Zhang W, Zhu Y (2021) Cross-modal image sentiment analysis via deep correlation of textual semantic. *Knowledge-Based Syst* 216:.. <https://doi.org/10.1016/j.knosys.2021.106803>
- [4] Lu J, Batra D, Parikh D, Lee S (2019) ViLBERT: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Adv Neural Inf Process Syst* 32:1–11
- [5] Xu J, Huang F, Zhang X, Wang S, Li C, Li Z, He Y (2019) Visual-textual sentiment classification with bi-directional multi-level attention networks. *Knowledge-Based Syst* 178:61–73. <https://doi.org/10.1016/j.knosys.2019.04.018>
- [6] Huang F, Zhang X, Zhao Z, Xu J, Li Z (2019) Image–text sentiment analysis via deep multimodal attentive fusion. *Knowledge-Based Syst* 167:26–37. <https://doi.org/10.1016/j.knosys.2019.01.019> (0123456789,-().vol
- [7] Li Z, Sun Q, Guo Q, Wu H, Deng L, Zhang Q, Zhang J, Zhang H, Chen Y (2021) Visual sentiment analysis based on image caption and adjective–noun–pair description. *Soft Computing.* [https://doi.org/10.1007/s00500-021-06530-6\(0123456789\(\).,-volV](https://doi.org/10.1007/s00500-021-06530-6(0123456789().,-volV)
- [8] Serra A, Carrara F, Tesconi M, Falchi F (2023) The Emotions of the crowd: Learning image sentiment from Tweets via cross-modal distillation. <https://doi.org/10.48550/arXiv.2304.14942>
- [9] Zhu T, Li L, Yang J, Zhao S, Liu H, Qian J (2022) Multimodal sentiment analysis with image-text interaction network. *IEEE Transactions on Multimedia* 25: 3375 – 3385
- [10] Yadav A, Vishwakarma D (2020) A deep multi-level attentive network for multimodal sentiment analysis. <https://doi.org/10.48550/arXiv.2012.08256>
- [11] Li Z, Xu B, Zhu C, Zhao T (2022) CLMLF: A contrastive learning and multi-layer fusion method for multimodal sentiment detection. <https://doi.org/10.48550/arXiv.2204.05515>
- [12] Peng C, Zhang C, Xue X, Gao J, Liang H, Niu Z (2022) Cross-modal complementary network with hierarchical fusion for multimodal sentiment classification. *TSINGHUA SCIENCE AND TECHNOLOGY* 27: 664:679
- [13] Xu N, Mao W (2017) MultiSentiNet: A Deep Semantic Network for Multimodal Sentiment Analysis. In *Proceeding of the 2017 ACM on Conference on Information and Knowledge Management (CIKM'17)* 2399-2402
- [14] Wen H, You S, Fu Y (2021) Cross-modal context-gated convolution for multi-modal sentiment analysis. *Pattern Recognit Lett* 146:252–259. <https://doi.org/10.1016/j.patrec.2021.03.025>
- [15] Yang B, Shao B, Wu L, Lin X (2022) Multimodal Sentiment Analysis with Unidirectional Modality Translation. *Neurocomputing.* <https://doi.org/10.1016/j.neucom.2021.09.041>
- [16] Rahman W, Hasan MK, Lee S, Zadeh A, Mao C, Morency L, Hoque E (2020) Integrating multimodal information in large pretrained transformers. *Proc Annu Meet Assoc Comput Linguist* 2359–2369. <https://doi.org/10.18653/v1/2020.acl-main.214>
- [17] Devlin J, Chang MW, Lee K, Toutanova K (2019) BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT 2019 - 2019 Conf North Am Chapter Assoc Comput Linguist Hum Lang Technol - Proc Conf* 1:4171–4186

- [18] Hazarika D, Zimmermann R, Poria S (2020) MISA: Modality-Invariant and -Specific Representations for Multimodal Sentiment Analysis. *MM 2020 - Proc 28th ACM Int Conf Multimed* 1122–1131. <https://doi.org/10.1145/3394171.3413678>
- [19] Mittal T, Bhattacharya U, Chandra R, Bera A, Manocha D (2020) M3ER: Multiplicative multimodal emotion recognition using facial, textual, and speech cues. *AAAI 2020 - 34th AAAI Conf Artif Intell* 1359–1367. <https://doi.org/10.1609/aaai.v34i02.5492>
- [20] Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani M, Minderer M, Heigold G, Gelly S, Uszkoreit J, Houlsby N (2020) An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. <https://doi.org/10.48550/arXiv.2010.11929>
- [21] Han K, Wang Y, Chen H, Chen X, Guo J, Liu Z, Tang Y, Xiao A (2023) A survey on visual transformer. <https://doi.org/10.1109/TPAMI.2022.3152247>
- [22] Reimers N, Gurevych I (2019) Sentence-BERT: Sentence embeddings using siamese BERT-networks. *EMNLP-IJCNLP 2019 - 2019 Conf Empir Methods Nat Lang Process 9th Int Jt Conf Nat Lang Process Proc Conf* 3982–3992. <https://doi.org/10.18653/v1/d19-1410>
- [23] Wang B, Kuo C (2020) SBERT-WK: A Sentence embedding method by dissecting BERT-based word models. <https://doi.org/10.48550/arXiv.2002.06652>
- [24] Tsai, Y, Bai S, Liang P, Kolter J, Morency L, Salakhutdinov R (2019) Multi-modal transformer for unaligned multimodal language sequences. *Proceedings of the Annual Meeting of the Association for Computational Linguistics* 6558–6569
- [25] Niu T, Zhu S, Pang L, Saddik A (2016) Sentiment analysis on multi-view social data. In *International Conference on Multimedia Modeling* 15–27
- [26] Keila D, Firooz H, Mohan A, Goswami V, Singh A, Ringshia P, Testuggine D (2020) The hateful memes challenge: Detecting hate speech in multimodal memes. *Annual Conference on Neural Information Processing Systems* 2611-2624
- [27] Onita D, Dinu L, Adriana B (2019) From image to text in sentiment analysis via regression and deep learning. *Proceedings of Recent Advances in Natural Language Processing* 862–868
- [28] Wang W, Wei F, Dong L, Bao H, Yang N, Zhau M (2020) MINILM: Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers
- [29] Cer D, Yang Y, Kong S, Hua N, Limtiaco N, John R, Constant N, Guajardo-Ce'spedes M, Yuan S, Tar C, Sung Y, Strope B, Kurzweil R (2018) Universal sentence encoder. <https://doi.org/10.48550/arXiv.1803.11175>, 2018
- [30] Le Q, Mikolov T (2014) Distributed Representations of Sentences and Documents. <https://doi.org/10.48550/arXiv.1405.4053>