ISSN: 2321-2152 **IJMECCE** International Journal of modern electronics and communication engineering

Charl

E-Mail editor.ijmece@gmail.com editor@ijmece.com

www.ijmece.com



IMSSO ALGORITHM WITH OPTIMIZED HDFS MODEL ON BIG DATA ANAMOLY DETECTION USING IHCNN

Mule Rama Krishna Reddy*, Dr. O Naga Raju

ABSTRACT:

This study presents a novel approach integrating the Social Spider Algorithm (SSA) and a Hybrid Convolutional Neural Network (CNN) to address memory optimization and reduce data redundancy in big data analysis, specifically applied to the KDD dataset. The KDD Cup dataset, a benchmark in intrusion detection systems, is known for its high dimensionality and complexity, which poses challenges for efficient analysis and model performance. Our proposed approach aims to mitigate these challenges by leveraging the SSA for its capacity to interpolate and optimize the search space while incorporating a Hybrid CNN architecture, merging the strengths of different CNN models to effectively extract features and reduce memory usage. The Social Spider Algorithm serves as a dynamic optimization tool inspired by the cooperative behaviour of social spiders, which fosters a collaborative search strategy for identifying critical data patterns. By applying the SSA in the feature selection process, the algorithm intelligently explores the dataset, reducing redundant and irrelevant features, thus optimizing memory consumption without compromising the data's integrity. Additionally, the proposed Hybrid CNN model harnesses the power of various CNN architectures, amalgamating their diverse capabilities to handle complex data structures more efficiently. This approach aids in feature extraction by identifying significant patterns within the reduced feature space generated by the SSA, ultimately enhancing the model's predictive performance. The experimental results demonstrate the efficacy of the proposed approach in terms of both memory optimization and predictive accuracy. The hybrid approach significantly reduces the memory footprint while maintaining or improving the predictive performance on the KDD dataset, showcasing the potential for practical application in big data analytics, particularly in intrusion detection systems and cybersecurity. This research contributes to advancing the field of big data analysis by offering a more memory-efficient and accurate model, thereby opening avenues for more effective data-driven decision-making in various domains.

Keywords: Social Spider Algorithm (SSA), Hybrid Convolutional Neural Network (CNN), Memory Optimization, Data Redundancy, Big Data Analysis, KDD Dataset, Intrusion Detection Systems, High Dimensionality, Feature Selection, Feature Extraction Predictive Performance, Memory Consumption, Experimental Results, Cybersecurity, Data Integrity

INTRODUCTION:

In the realm of big data analytics, diverse machine learning and deep learning algorithms play a vital role in extracting meaningful insights from complex datasets. Besides the innovative approach combining the Social Spider Algorithm (SSA) with an interpolative hybrid Convolutional Neural Network (CNN) for the KDD dataset, several other methodologies contribute significantly to optimizing memory usage and handling redundant features across various domains. Gradient Boosting Machines (GBM) and Random Forest are widely acknowledged for their efficacy in handling large volumes of data. These algorithms are proficient in memory management by constructing ensembles of decision trees, effectively reducing redundancy in features and providing robust predictive modeling

* Research Scholar, Dept of Computer Science and Engineering, Acharya Nagarjuna University, Guntur, Andhra Pradesh. Associate Professor, Head - Dept of Computer Science, Government Degree College, Dornala, Andhra Pradesh



ISSN2321-2152 www.ijmece .com Vol 14, Issue. 4 Nov2023

Their adaptability to various data types and feature engineering makes them advantageous in scenarios where memory optimization and feature extraction are paramount. Support Vector Machines (SVM) offer another avenue in big data analysis by creating decision boundaries for classification or regression tasks. Their ability to reduce memory usage by selecting a subset of training instances as support vectors while ensuring robustness in highdimensional spaces makes them suitable for complex datasets like KDD. Additionally, their kernel trick facilitates efficient feature extraction by transforming higher dimensions for data into improved classification accuracy. In the domain of deep learning, Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) networks provide sophisticated memory handling and feature extraction capabilities. Particularly in sequential data analysis within intrusion detection systems, these models excel at learning temporal dependencies, reducing redundancy in sequential features, and effectively capturing patterns within the data. Furthermore, Autoencoders, a type of neural network, contribute significantly by learning data representations in an unsupervised manner. They aid in memory optimization by extracting essential features through dimensionality reduction, thus mitigating redundant information and enabling improved efficiency in subsequent analysis. Ensemble learning methods such as Stacking and Bagging techniques combine multiple models to enhance predictive accuracy while effectively managing memory usage and reducing redundant features. These approaches optimize data analysis by leveraging the strengths of various algorithms, improving overall model performance.

Each of these algorithms and methodologies presents unique strengths in addressing memory optimization and redundant features in big data analysis. The synergy and integration of these diverse approaches can offer comprehensive solutions to the challenges faced in various domains, fostering more effective data-driven decision-making and analysis across industries like cybersecurity, finance, healthcare, and beyond. The exploration and fusion of these methodologies with the SSA and Hybrid CNN approach represent a progressive direction in advancing big data analysis methodologies, aiming to unlock new possibilities for efficient data handling and predictive modelling which are mentioned in this paper to implicate the need of Learning algorithms on big data.

Challenges in Big Data Analysis:

The rapid proliferation of data across various industries has given rise to an urgent need for advanced analytics capable of extracting meaningful insights from this ever-expanding wealth of information. The KDD Cup dataset, renowned for its application in intrusion detection systems, serves as an exemplar of complex, high-dimensional data. However, this richness presents inherent challenges, notably the intricate nature of the dataset, which demands sophisticated approaches for effective analysis. The dataset's high dimensionality complicates the task of efficient memory usage, making it imperative to optimize the computational resources for streamlined analysis. Moreover, the presence of redundant features within the dataset adds to the complexity, hampering the accuracy and efficiency of modelling. These challenges form the bedrock of this study, motivating the exploration of innovative methodologies that can effectively address memory optimization and feature redundancy within the context of the KDD dataset.

The Promise of the Social Spider Algorithm and Interpolative Hybrid CNN:

The Social Spider Algorithm (SSA) emerges as a promising tool in addressing the intricacies of data redundancy and memory optimization. Drawing inspiration from the cooperative behaviour of social spiders, this algorithmic framework operates as a dynamic optimization mechanism capable of identifying and isolating critical data patterns. By intelligently exploring the dataset, SSA facilitates the identification and reduction of redundant features, thereby enhancing memory efficiency without compromising the dataset's integrity. The SSA's adaptive nature aligns with the data exploration process, aiding in the selection of relevant features, a fundamental step towards memory optimization. Simultaneously, the Interpolative Hybrid CNN approach represents an innovative strategy that leverages the prowess of various Convolutional Neural Network architectures. This approach aims to



amalgamate the strengths of different CNN models, utilizing their varied capabilities to extract essential features while concurrently reducing memory usage. By consolidating the benefits of diverse CNN architectures, this hybrid model seeks to enhance feature extraction within the reduced feature space identified by the SSA. The synergy between these two methodologies is envisioned to not only improve predictive performance but also significantly reduce the memory footprint, a critical advancement in the realm of big data analysis, especially within intrusion detection systems and cybersecurity.

Potential Impact and Future Applications:

The potential impact of this research is far-reaching, transcending the immediate domain of intrusion detection systems and cybersecurity. By developing a memory-efficient and accurate model that excels in handling high-dimensional and complex datasets, this research lays the foundation for broader applications across various industries. The implications extend to healthcare, finance, marketing, and other sectors heavily reliant on data analytics. A successful integration of the SSA and the Hybrid CNN model in mitigating memory constraints and enhancing predictive accuracy could revolutionize decisionmaking processes reliant on big data analysis, empowering stakeholders with more precise and actionable insights. This research, therefore, not only offers solutions for the challenges posed by the KDD dataset but also opens doors to a new era of more effective, data-driven decision-making in a multitude of industries.

Problem statement

The challenge within big data analysis lies in the KDD dataset's intricate nature, characterized by high dimensionality and data redundancy, impeding efficient memory usage and accurate predictive modeling, particularly in intrusion detection systems. This study aims to tackle these obstacles by integrating the Social Spider Algorithm (SSA) with an interpolative hybrid Convolutional Neural Network (CNN), with a specific focus on reducing memory consumption and eliminating redundant features in the KDD dataset. The overarching goal is to develop an interpolative model that optimizes memory utilization and feature extraction.

significantly enhancing the accuracy and efficiency of predictive modeling within the complex landscape of intrusion detection systems.

Objectives:

- 1. To implement the Social Spider Algorithm (SSA) in conjunction with an interpolative hybrid Convolutional Neural Network (CNN) for the purpose of identifying and eliminating redundant features within the KDD dataset. The primary aim is to optimize memory utilization by intelligently reducing the dataset's dimensionality, preserving critical data patterns, and enhancing the overall efficiency of the model in intrusion detection systems.
- 2. To develop an interpolative model that enhances feature extraction within the KDD dataset through the collaborative application of the SSA and a hybrid CNN approach. This objective targets the reduction of memory footprint while significantly improving predictive accuracy. By synergizing the strengths of the SSA and diverse CNN architectures, the objective is to create a model that not only reduces redundant features but also maximizes the extraction of essential patterns, thereby advancing the efficacy of intrusion detection and cybersecurity systems.

Overview:

This paper is structured as follows: following this introduction, the subsequent sections will delve into the theoretical framework underlying the Social Spider Algorithm and the Interpolative Hybrid CNN Subsequently, approach. the experimental methodology, results, and discussions will be presented, culminating in а comprehensive conclusion that outlines the implications and future directions of this research.

LITERATURE SURVEY:

The research community has extensively utilised artificial intelligence, namely machine learning, to transform a wide range of disparate and heterogeneous data sources into reliable facts and



knowledge. This has resulted in the development of advanced capabilities for accurately identifying patterns. Nevertheless, the utilisation of machine learning techniques on extensive and intricate datasets incurs significant computational costs, necessitating substantial allocation of logical and physical resources, including data file storage, central processing unit (CPU), and memory. The significance of a sophisticated platform for efficient big data analytics has increased in recent times due to the exponential growth in the volume of data generated on a daily basis, which now exceeds quintillion bytes. Apache Spark MLlib is widely recognised as a leading platform for the analysis of large-scale datasets. It provides a comprehensive range of advanced functionalities for various machine learning encompassing regression, classification, tasks, dimension reduction, clustering, and rule extraction. This paper examines the computational aspects of the Apache Spark MLlib 2.0, which is an open-source machine learning framework that is distributed, scalable, and platform independent. In this study, a series of practical machine learning experiments are conducted to investigate the qualitative and quantitative characteristics of the platform. In addition, we emphasise the prevailing patterns in research on machine learning in the context of big data and offer valuable perspectives for prospective investigations.

The dataset has three distinct types of comparison matrices. Three distinct datasets have been made available for the purpose of analysing the tangible effects of machine learning within the realm of big data analytics. The datasets provide an overview of various machine learning methodologies, the performance of machine learning algorithms, and a comparative assessment of big data technology. We have conducted an analysis on various machinelearning methodologies and subsequently generated these datasets. The utilisation of these datasets is commendable for constructing solutions using emerging technologies in the field of bug data. The selection of division criteria for decision trees is based on many quality indicators, necessitating the consideration of the entire dataset for each branching node. The utilisation of decision trees in large-scale data applications is significantly hindered. Support Vector Machines (SVM) demonstrate exceptional performance when used to moderate-sized datasets containing information collections. There are inherent obstacles that hinder the implementation of largescale information applications. Profound learning is well-suited to address challenges associated with the scale and diversity of large datasets. However, it is important to note that there are certain restrictions associated with handling big amounts of information, since it necessitates a significant amount of time for processing.

In recent years, there has been an increasing focus on the application of big data and Internet of Things (IoT) technologies. The primary objective of the researchers was the development of big data analytics solutions through the utilisation of machine learning models. The utilisation of machine learning has been increasingly prevalent in this domain owing to its capacity to discern latent characteristics and patterns inside intricate datasets. This study employed the Big Data IoT Framework to conduct an analysis of weather data in a specific use case. The implementation of weather clustering and sensor anomaly detection was conducted using a publically accessible dataset. The implementation details for each tier of the architecture (collection, ETL, data processing, learning, and decision) were supplied for this specific use case. The learning model selected for implementation in the library is k-means clustering, which is based on the Scikit-Learn framework. The findings of the data analysis demonstrate the feasibility of extracting significant insights from a dataset of considerable complexity through the utilisation of our framework.

This paper introduces the Real-time Machine Learning Competition on Data Streams, which was held as part of the IEEE Big Data 2019 conference and was referred to as the BigData Cup Challenge. Data streams, particularly those generated by sensors, have garnered significant attention from both scholars and companies, and are currently the subject of extensive investigation within the field of data science. Companies operating in the telecommunication and energy sectors are actively seeking to leverage these data in order to gain timely and valuable insights pertaining to their respective services and equipment. To effectively derive



important insights from data streams, it is imperative to possess the capability to analyse the incoming data in real-time and provide relevant forecasts. Fast incremental learners are utilised for this particular objective. There is currently an established community that is actively engaged in organising a diverse range of competitions focused on machine learning tasks specifically designed for batch learners. The objective of our study was to implement a similar methodology aimed at including the entire community in addressing critical issues in the field of data stream mining. A novel data science competition was conducted, wherein a real-time prediction scenario was employed. This competition utilised a unique competition platform specifically designed for data streams. The instances for prediction were made available in real-time, and the corresponding forecasts were required to be submitted in real-time as well. Based on the available information, it can be inferred that this particular data science competition is the inaugural instance of a real-time format being employed. The objective of the competition entailed forecasting network activity, with the dataset being furnished by one of our affiliated corporate entities.

In contemporary times, a substantial amount of data is being generated on a daily basis, measured in gigabytes. This data exhibits various characteristics, including but not limited to high velocity, extensive volume, inherent uncertainty, non-stationary nature, and real-time availability. The study of huge data is not feasible using conventional machine learning techniques, as previously stated. Moreover, conventional storage and processing methods are inadequate in meeting the specified requirements. This study examines the issues associated with utilising classical machine learning techniques (MLT) for Big Data Analytics and proposes potential ways to address these challenges. According to the findings of our survey, several potential solutions have been identified to address the issues inherent in big data analytics. These solutions include parallel processing, dimensionality reduction techniques, the utilisation of GPUs, the implementation of map reduce jobs, the application of deep learning methodologies, as well as online learning and incremental learning approaches.

In contrast to intelligent analysis applications in traditional small-scale data scenarios, intelligent analysis applications in big data scenarios present a distinct challenge. They no longer revolve around a single AI algorithm model, but rather involve the integration of big data, big models, and big computing. The current data fusion technique exhibits several limitations, including elevated network energy consumption and a reduced lifespan of network nodes subsequent to data fusion. Simultaneously considering algorithm model design, big data processing, and efficient distributed parallel computing is imperative. This integration introduces a multitude of novel challenges and issues to the study of fundamental theoretical approaches and pivotal technologies in the realm of intelligent analysis of big data. The methodology employed in this study involves the utilisation of machine learning techniques, wherein an algorithm is implemented to enable autonomous data analysis through deep mining processes. It has the capacity to effectively complete tasks that are not feasible to be executed manually. To achieve comprehensive data analysis, it is vital to employ complicated analysis techniques that are grounded in machine learning and data mining. Additionally, the utilisation of robust processing power and ample storage capacity is essential in order to effectively handle and process vast amounts of data. Based on the findings of this study, it can be inferred that the method proposed in this paper exhibits superior temporal performance compared to the conventional approach.

The global community is confronted with the challenges posed by dynamic weather patterns and their associated consequences. In order to mitigate these adverse effects to a certain extent, numerous techniques and algorithms have been developed. These tools enable us to forecast weather conditions based on historical data, such as temperature, dew point, humidity, air pressure, and wind direction, thereby providing valuable information for preparedness purposes. During the analysis of historical data from recent years, we incorporated the proposed scheme or techniques that aim to demonstrate that the utilisation of machine learning paradigm enables us to investigate the provided dataset and extract valuable information.



Consequently, in order to comprehend the fluctuating patterns of climatic conditions, a predictive model is also advocated. In this paper or study, we will examine advanced statistical linear regression and support vector machine techniques in machine learning for predicting weather forecasts. These strategies involve utilising consistent datasets to make accurate predictions. The proposed scheme aims to incorporate an enhanced algorithm that generates approximate and proximate climate forecasts for the upcoming five days. The final results are determined through the utilisation of mathematical and statistical decision trees, as well as the application of a confusion matrix, in order to achieve more precise and accurate predictions. This approach leverages Big Data for forecasting purposes.

Individuals from many geographical locations have the opportunity to articulate their perspectives and viewpoints via a multitude of online social media platforms. Individuals engage in the everyday usage of online social media platforms as a means of interpersonal communication and to be abreast of contemporary occurrences. Every day, Twitter receives a substantial volume of tweets encompassing a diverse array of topics. Twitter is a highly recognised and extensively used digital social media network. The processes of feature extraction and trend identification can be effectively achieved by leveraging machine learning techniques. In order to effectively extract valuable insights from the continuous influx of data generated by Twitter, it is imperative to employ specialised tools and procedures tailored for handling substantial data volumes. This study primarily centres on the identification of hashtags and the determination of the industry with the greatest share of voice. This study aims to gather real-time data from the social media platform Twitter through the use of Apache Spark. Subsequently, the classification of each tweet is conducted through the utilisation of machine learning techniques offered by the Apache Spark machine learning library. In order to evaluate the performance of the model, a convolutional neural network (CNN) and logistic regression (LR) are employed. The convolutional neural network (CNN) approach demonstrated superior performance

compared to the logistic regression strategy, with an average accuracy of about 95% and an F1 score of 0.60. The current values for both accuracy and the F1 score are 0.59. Based on the results obtained, it has been observed that the use of the Apache Spark framework for big data enables a significantly faster evaluation of real-time tweets compared to the traditional execution environment. The findings indicate that the use of the Apache Spark tool for big data enables significantly faster evaluation of realtime tweets compared to conventional execution environments.

Healthcare plays a crucial role within the contemporary medical landscape, particularly in the context of the digital era. In the context of sickness prediction and other healthcare-related tasks, it is imperative for a healthcare system to thoroughly analyse extensive volumes of patient data. An intelligent system would possess the capability to examine several aspects of a patient's life, including their social interactions, medical background, and other characteristics related to their lifestyle, in order to predict the probability of encountering a health issue. The Health Recommender System (HRS) is gaining increasing importance as a means of delivering healthcare services. Health-intelligent systems have become essential elements in the decision-making process of healthcare delivery within this particular context. The main objective of their work revolves around ensuring the consistent provision of information that is of superior quality, dependable, genuine, and confidential, thereby enabling its optimal use. The health recommender system plays a vital role in generating outcomes such as suggesting diagnoses, health insurance options, treatment approaches based on clinical pathways, and alternative medications, all tailored to the patient's health profile. This is particularly significant as an increasing number of individuals are turning to social networks as a means to acquire health-related knowledge. To optimise healthcare efficiency and cost-effectiveness, recent research has concentrated on leveraging extensive medical data by integrating multimodal data from various sources. In the healthcare industry, the utilisation of big data analytics in conjunction with recommender systems holds significant importance when making decisions



pertaining to a patient's health. This article proposes the use of a LeNET Convolutional Neural Network (CNN) to explore the integration of big data analysis into the creation of effective health recommendation systems. It also demonstrates the potential advantages for the healthcare industry in transitioning from a standardised model to a personalised approach within the realm of telemedicine. The proposed approach demonstrates superior performance compared to other approaches by considering both the root squared mean error (RSME) and the average absolute error (AAE).

Colour constancy refers to the perceptual phenomenon in which individuals are able to accurately perceive and recognise the colours of objects, regardless of variations in the properties of the illuminating light source. The objective of computational colour constancy is to estimate the illuminant and afterwards use this knowledge to rectify the image and present it as it would appear under a standard illuminant. The deep learning approach has emerged as one of the most effective techniques for estimating illumination in various scenarios. This method normally requires a dataset of photos that have been annotated with the corresponding scene illumination. While drawing parallels between the human visual system and machine learning algorithms is common, it is important to note that the former has not been exposed to definitive and accurate information on illuminants throughout its evolutionary process. Alternatively, it is postulated that the emergence of colour constancy can be attributed to its facilitation of various essential functions, including the autonomous recognition of fruits, objects, and animals irrespective of the prevailing illumination conditions. The rapid advancement of artificial intelligence and the consequent enhancement of individuals' well-being have led to significant advances in the field of picture identification, particularly in recent years. This study investigates the problem of object detection in low illumination conditions and employs deep learning techniques for the purpose of image detection and analysis. In environments with limited lighting conditions, the detected items are subjected to a comparison process with a large dataset. This comparison aims to identify

and validate the objects that exhibit the highest degree of similarity. Additionally, this comparison evaluates the learning model's accuracy for object recognition in difficult lighting conditions.

The utilisation of mobile health has evolved as a viable and pragmatic solution for the treatment and management of individuals' health issues. However, the majority of mobile health data consists of observational data obtained by sensors, posing challenges in analysing the causal relationship between provided interventions using conventional regression techniques. This study provides a comprehensive overview of deep learning models that have the potential to accurately evaluate the causal effect of unprocessed mobile health data. The aforementioned models possess the ability to effectively process multivariate time series data for the purpose of estimating the unbiased causal effect, particularly when presented with a series of treatments.

Performance measurement systems play a crucial role in the administration of organisations since they facilitate the transformation of raw data into meaningful information that can be utilised by decision-makers. In the past few decades, there has been a significant surge in the volume of data and information produced and disseminated, presenting novel prospects and complexities for these systems. In light of the given situation, the objective of this essay is to examine the utilisation of big data analytics within performance measurement systems in order to elucidate the relationship between the two. Moreover, the objective is to discern patterns and potential avenues for future scholarly investigation. In order to accomplish this objective, a scientific mapping was conducted using bibliometric analytic techniques. The primary findings of the study indicate a notable rise in the utilisation of big data analytics within performance measurement systems (PMS) in recent years, without due consideration for the inherent characteristics of such systems. The integration of artificial intelligence technologies, specifically machine learning and deep learning, has the potential to enhance the field, opening up avenues



for empirical research using unstructured data and applications within the context of Industry 4.0.

Machine learning techniques, particularly natural language processing (NLP), are of significant importance in the context of using social media data for governmental purposes in many countries worldwide. The analysis of social media posts and tweets can provide insights on the prevailing thinking of individuals, a crucial aspect for any governing body worldwide. The primary aim of this research is to do sentiment analysis in order to extract the prevailing sentiment among individuals in relation to the ongoing conflict between Russia and Ukraine. This will be accomplished through the use of machine learning methodologies. The objective is to conduct an analysis and draw inferences regarding the potential reactions of countries in response to the economic effects, taking into account the sentiment of their respective citizens. The implementation process commences with the acquisition of data from social media platforms, specifically Twitter and Reddit. This is achieved through the use of Snscraper, a web scraping tool, and the PRAW (Python Reddit API Wrapper) library. Suitable text summarising techniques are employed to accommodate the larger posts on Reddit. The BERT transformer model is used to conduct sentiment analysis on social media data. The non-English posts undergo translation into with a deep learning-based short-term traffic forecasting (GWODL-STTF) model in the context of a smart city setting. The GWODL-STTF technique focuses on forecasting traffic flow in smart cities. The GWODL-STTF approach encompasses two primary steps. In the early phase, the GWODL-STTF methodology utilised a gated recurrent unit-neural network (GRU-NN) model for the purpose of predicting traffic flow. In the subsequent phase, the GWODLSTTF approach employs the Grey Wolf Optimisation (GWO) algorithm as a hyperparameter optimizer. The performance of the GWODL-STTF technique may be evaluated using many metrics in simulation experiments. The minimum mean squared error (MSE) value of 105.627 obtained from the

results shows that the GWODL-STTF method

performs better than recent techniques.

English through the use of neural machine translation. In addition, sentiment analysis is conducted at several levels of detail, including the identification of specific locations and individuals through the application of named entity recognition algorithms. In this study, a comprehensive examination is conducted to compare the emotions of various countries throughout the world with their respective levels of dependence on Russian oil.

The Intelligent Transportation System (ITS) is a groundbreaking technology within the realm of smart cities that serves to mitigate traffic congestion and enhance traffic conditions. Information Technology Services (ITS) offers real-time analysis and highly efficient traffic control through the utilisation of big data and communication technology. Traffic flow prediction (TFP) has emerged as a crucial element in the management of smart cities, serving as a means to forecast forthcoming traffic conditions on transport networks based on historical data. Machine learning (ML) and neural network (NN) methodologies have demonstrated significant use in addressing real-time challenges due to their ability to effectively handle dynamic data over extended periods. Deep learning (DL) is a subset of machine learning (ML) techniques that demonstrate high efficacy in tasks related to prediction and data classification. This article presents the development of a Grey Wolf optimizer

Medical image classifiers serve a key role in both medical services and educational endeavours. However, the conventional technique reached its maximum level of performance. Furthermore, the utilisation of these traits necessitates a significantly greater amount of time and effort for extraction and selection. The Deep Neural Network (DNN) is an emerging machine learning (ML) technique that has demonstrated its potential for many classification problems. The convolutional neural network (CNN) has been found to yield optimal results in several image classification tasks. However, the acquisition of medical picture databases can pose challenges due to the need for specialised expertise in categorization. This research paper presents the development of a novel hyperparameter-tuned deep learning model for healthcare monitoring systems (HPTDLM-HMS) within the context of a big data environment. The HPTDLM-HMS technique discussed in this study



focuses on the analysis of medical pictures within the context of decision-making. The HPTDLM-HMS technique is initially implemented by utilising the EfficientNet model to extract features. The hyperparameters of the model are optimised using the Manta Ray Foraging Optimisation (MRFO) algorithm. Finally, the categorization of medical images is conducted using the Long Short-Term Memory (LSTM) technique. Hadoop MapReduce is employed for the purpose of managing large volumes of data. The outcome evaluation of the HPTDLM-HMS approach is assessed using a dataset consisting of medical imaging data. The full examination of the HPTDLM-HMS approach has demonstrated a recall value of 87.46%, which surpasses the performance of alternative models and underscores its potential prospects.

In contemporary times, there has been a substantial proliferation of data, leading to a progressive transformation in the importance attributed to data security and data analysis techniques within the realm of big data." An intrusion detection system (IDS) is a mechanism that examines and monitors data in order to identify any unauthorised access or incursion into a system or network. The substantial magnitude, diversity, and rapid velocity of data generated within the network necessitate a sophisticated data analysis methodology to effectively identify and mitigate assaults. Big data systems can be employed in intrusion detection systems (IDS) to facilitate the management of large volumes of data, enabling accurate and efficient data analysis methodologies. This research paper presents a novel approach called Intrusion Detection Approach utilising Hierarchical Learning-based Butterfly Optimisation Deep Algorithm (ID-HDLBOA) in the context of a big data platform. The technique known as ID-HDLBOA integrates the principles of deep learning (DL) with the process of hyperparameter tuning. The ID-HDLBOA technique incorporates a hierarchical LSTM model for the purpose of intrusion detection. The BOA method is employed as a hyperparameter tuning strategy for the LSTM model, leading to enhanced detection efficiency. The ID-HDLBOA technique is experimentally validated using a benchmark incursion dataset, yielding a model accuracy of 98%. A series of comprehensive experiments were conducted, and the results consistently highlighted the superior performance of the ID-HDLBOA algorithm.

This paper presents a proposed recommendation system, named Sentiment Analysis and Matrix Factorization (SAMF), which aims to address the challenges of data sparsity and credibility in collaborative filtering. SAMF leverages topic modelling and deep learning techniques to effectively extract implicit information from reviews. By enhancing the rating matrix, SAMF assists in improving the recommendation process. The generation of user topic distribution and item topic distribution is accomplished by applying Latent Dirichlet Allocation (LDA) to reviews, which include both user reviews and item reviews. The user feature matrix and item feature matrix are generated by utilising topic probability. Furthermore, the integration of the user feature matrix and item feature matrix results in the creation of a user-item preference matrix. In addition, the process involves the integration of the user-item preference matrix and the original rating matrix, resulting in the creation of the user-item rating matrix. In addition, the utilisation of BERT (Bidirectional Encoder Representation from Transformers) is employed to quantify the sentiment information encompassed within the reviews. This sentiment information is subsequently integrated with user-item rating matrix, facilitating the the modification and updating of said matrix. Subsequently, the revised user-item rating matrix is employed to facilitate the prediction of ratings and generate Top-N recommendations. The experimental results obtained by analysing Amazon datasets provide evidence that the proposed SAMF algorithm outperforms existing conventional algorithms in terms of suggestion performance.

In recent years, there has been a significant utilisation of deep learning (DL) in academic research pertaining to the interpretation of 12-lead electrocardiogram (ECG) data. Nevertheless, there is a lack of clarity regarding the validity of the explicit or implicit assertions regarding the superiority of deep learning (DL) over classical feature engineering (FE) approaches that rely on domain expertise. Furthermore, there is a lack of clarity on the potential enhancement of performance through the integration



of deep learning (DL) and feature engineering (FE) in comparison to utilising a single modality. Methods: In order to fill the existing research gaps and align with recent significant experiments, we conducted a reexamination of three specific tasks: the diagnosis of cardiac arrhythmia using a multiclass-multilabel classification approach; the prediction of atrial fibrillation risk using a binary classification approach; and the estimation of age using a regression approach. For the purpose of our study, we utilised a comprehensive dataset consisting of 2.3 million 12-lead electrocardiogram (ECG) recordings. These recordings were employed to train various models for each specific task. Specifically, we developed three models: i) a random forest model that utilised feature extraction (FE) as input; ii) an end-to-end deep learning (DL) model; and iii) a merged model that combined both FE and DL approaches. The findings indicate that the performance of feature engineering (FE) was similar to that of deep learning (DL) in the two classification tasks. However, FE required much less data compared to DL. The deep learning (DL) approach demonstrated superior performance compared to the traditional feature engineering (FE) method for the regression problem. In all tasks, the integration of feature engineering (FE) with deep learning (DL) did not yield any performance improvement compared to the use of DL alone. The aforementioned conclusions were validated using the supplementary PTB-XL dataset. In conclusion, our findings indicate that deep learning (DL) did not demonstrate a substantial enhancement over feature engineering (FE) in the context of traditional 12-lead electrocardiogram (ECG)-based diagnosis tasks. However, it did exhibit a notable improvement in the atypical regression challenge. Furthermore, our findings indicate that the integration of feature engineering (FE) with deep learning (DL) did not yield superior results compared to DL in isolation. This observation shows that the features extracted using FE were duplicative of the features acquired through DL. The significance of our findings lies in the provision of crucial advice pertaining to the selection of a suitable machine learning technique and data regime for a specific task, with a focus on 12-lead ECG analysis. When considering the objective of maximising performance, deep learning (DL) is the preferred approach when dealing with unconventional tasks and large datasets. In cases where the work at hand is of a classical nature and/or there is limited availability of data, employing a feature engineering (FE) technique may prove to be a more suitable option.

Accurate identification of heart disease can have lifesaving implications, while an erroneous diagnosis can have fatal consequences. The UCI dataset on heart disease in machine learning serves as a platform for evaluating and comparing the outcomes and analyses of different machine learning methodologies, encompassing deep learning techniques. The research was conducted using a dataset of 13 major characteristics. The datasets are processed using support vector machines and logistic regression techniques, with the latter demonstrating superior accuracy in predicting coronary disease. Python programming is employed for the purpose of processing datasets. Several research initiatives have employed machine learning techniques to enhance the efficiencv of the healthcare industry. Conventional machine learning methods were employed in our study to elucidate the relationships between the various variables included in the dataset. These findings were subsequently utilised to efficiently predict the risks of heart infections. The utilisation of the accuracy and confusion matrix has vielded positive benefits. In order to optimise the the dataset incorporates outcomes. specific extraneous attributes, which are addressed through the utilisation of isolation logistic regression and support vector machine (SVM) classification techniques.

EXISTING SSO APPORACH ON HDFS

The application of the SSO approach to HDFS for the KDD dataset for memory optimization involves a novel integration of machine learning algorithms in a distributed file system environment. This theoretical approach aims to leverage the principles of SSO, inspired by the cooperative behavior of social spiders, to optimize memory utilization while handling the high-dimensional and complex nature of the KDD dataset within the HDFS framework.

Principal Component Analysis (PCA) PCA, a dimensionality reduction technique, can be employed



to reduce the number of features within the KDD dataset, minimizing the memory footprint in HDFS. By transforming the data into a lower-dimensional space while retaining essential information, PCA can help alleviate memory strain without significant loss of predictive power.

Sampling Methods Utilizing sampling methods such as Random Sampling or Stratified Sampling can enable machine learning algorithms to work on smaller representative subsets of the KDD dataset within HDFS. This approach assists in reducing memory consumption while maintaining the statistical properties of the original dataset.

Compression Techniques Applying model compression techniques like Quantization or Pruning within the machine learning algorithms can significantly reduce the memory requirements in HDFS. Compact data structures and optimized storage mechanisms for models, particularly in deep learning models applied to the KDD dataset, can help alleviate memory strain.

Distributed Computing Leveraging distributed computing frameworks such as Apache Spark over HDFS enables parallel processing, allowing machine learning algorithms to distribute memory and computational load across multiple nodes. This aids in handling the large volume of the KDD dataset more efficiently, reducing the burden on individual systems.

Online Learning Algorithms Utilizing online learning algorithms within HDFS for the KDD dataset can dynamically update models based on incoming data, thus requiring less memory for storage and processing. Algorithms like Online Gradient Descent or Online Random Forests adapt to changing data, offering memory efficiency in evolving datasets.

Memory-Efficient Models Selecting memoryefficient models designed for large datasets within HDFS, such as LightGBM or XGBoost, can significantly reduce memory usage. These algorithms are optimized for both memory and computational efficiency, making them well-suited for handling the complexities of the KDD dataset. **Data Streaming Techniques** Implementing streaming algorithms in HDFS for real-time data processing of the KDD dataset can help manage memory usage. Streaming methods like Apache Flink or Apache Kafka offer continuous and manageable data processing, aiding in memory optimization.

Dynamic Resource Allocation Utilizing dynamic resource allocation strategies within HDFS for different machine learning algorithms working on the KDD dataset can optimize memory usage. Efficient resource management, such as Hadoop YARN, dynamically allocates resources as per the changing demands of algorithms, ensuring optimal memory utilization.

Continuous Optimization and Research The application of the SSO approach with machine learning algorithms for memory optimization in HDFS and the KDD dataset should be an ongoing process. Continuous exploration and research into new algorithms, distributed computing techniques, and innovative methodologies are essential for further advancements in memory optimization and efficient handling of large-scale datasets within the Hadoop environment.

DESIGN METHODOLOGY

The SSO approach (Social Spider Optimization) is an algorithm inspired by the collaborative behavior of social spiders, applied within the domain of machine learning and data analysis, particularly for the KDD dataset, a benchmark in intrusion detection systems. In this context, an analysis of the existing SSO approach alongside various machine learning algorithms is essential for problem identification within the KDD dataset.

Overview of SSO in Machine Learning Algorithms:

The SSO algorithm operates by mimicking the collaborative behavior of social spiders, creating an intelligent optimization strategy. It focuses on leveraging the algorithm's ability to explore complex datasets and identify critical data patterns. This collaborative optimization tool has shown promise in the realm of intrusion detection systems, like the



KDD dataset, known for its high dimensionality and complexity.

The existing Social Spider Optimization (SSO) approach integrated with various machine learning algorithms, including Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Random Forest Classifier (RFC), Decision Trees (DT), and other boosting algorithms, has proven beneficial in addressing anomaly detection challenges within the KDD dataset, a benchmark in intrusion detection systems.

- 1. **SVM and SSO Integration:** The amalgamation of SSO with SVM leverages SSO's collaborative optimization, enhancing SVM's performance in handling the complex, high-dimensional nature of the KDD dataset. SSO assists in identifying critical data patterns for improved anomaly detection.
- 2. KNN and SSO Collaboration: Integrating SSO with KNN aids in refining the dataset's features and improving proximity-based classifications. SSO dynamically optimizes the KNN's clustering and classification processes, enhancing its efficiency in identifying network intrusions.
- 3. **RFC and SSO Synergy:** The fusion of SSO with RFC enables the optimization of ensemble learning within the RFC model. SSO contributes to the decision-making process by refining the ensemble learning mechanism, resulting in more accurate and robust detection of anomalies.
- 4. **DT and SSO Complementarity:** SSO's intelligent optimization strategy complements DT by aiding in the decision-making process at each node of the tree. This collaboration optimizes feature selection and enhances the interpretability of the decision trees within the KDD dataset.

FLOW DIAGRAM

- 5. Boosting Algorithms and SSO Enhancement: Other boosting algorithms, such as AdaBoost or Gradient Boosting, integrated with SSO, exhibit improved performance in identifying and flagging anomalies. SSO assists in refining the weak learners and ensemble models, leading to more accurate detection results.
- 6. **Model Evaluation and Parameter Tuning:** The combined SSO with each algorithm undergoes comprehensive evaluation and fine-tuning of parameters to optimize anomaly detection accuracy on the KDD dataset. This process enhances the robustness and efficiency of the models.
- 7. Feature Selection and Dimensionality Reduction: SSO aids in feature selection and dimensionality reduction, significantly reducing redundant features within the KDD dataset. This facilitates a more streamlined and efficient learning process for all integrated algorithms.
- 8. Scalability and Real-world Applicability: The collective utilization of SSO with a range of machine learning algorithms ensures scalability and real-world applicability. This integrated approach shows promise in enhancing security and reliability in intrusion detection systems by handling the challenges of the intricate KDD dataset.

This collaborative integration of the Social Spider Optimization approach with various machine learning algorithms demonstrates the potential to enhance anomaly detection capabilities within the KDD dataset. It underscores the importance of leveraging the collective strengths of SSO and diverse algorithms to address the complexities and challenges inherent in intrusion detection systems.





PROPOSED MODEL:

The concept of an interpolative model leveraging Social Spider Optimization (SSO) principles for memory optimization within Hadoop Distributed File System (HDFS) involves utilizing an adaptive approach inspired by the collaborative behavior of social spiders to efficiently manage memory utilization. While there might not be explicit equations for this concept due to its heuristic nature, the essence involves interpolation and optimization techniques applied in a distributed computing framework.

Interpolative Memory Optimization:

The application of interpolation within the context of memory optimization aims to efficiently utilize memory resources within the distributed environment of HDFS. It involves the intelligent 'interpolation' between memory requirements and data access patterns to optimize the allocation and utilization of memory resources, mimicking the cooperative behavior observed in social spiders.

Social Spider Optimization (SSO) Adaptation:

SSO principles, derived from the collaboration patterns of social spiders, could be adapted to an interpolative model in the HDFS environment. This adaptation might involve heuristic rules inspired by social spiders' collective movement, such as exploring, clustering, or adaptively searching for memory allocation patterns, taking cues from the spiders' cooperative behavior to optimize memory usage.

Abstract Representation:

While not a direct mathematical equation, a conceptual representation could involve adaptive interpolation of memory utilization within HDFS:

Moptimized=f(*Mcurrent*,*Dpatterns*)

Here, *Moptimized* represents the optimized memory allocation, a function $\oint f$ adaptively interpolates between the current memory allocation (*Mcurrent*) and data access patterns (*Dpatterns*) within HDFS. This adaptive function could be guided by heuristics inspired by SSO to dynamically adjust memory utilization based on evolving data patterns.

Implementation and Heuristics:

The actual application of the concept would involve developing heuristics inspired by the collaborative behavior observed in social spiders. These heuristics might guide the interpolation between current memory utilization and changing data access patterns within HDFS, dynamically optimizing memory allocation.



Overall, an interpolative model within HDFS leveraging SSO concepts would require a sophisticated adaptation of the heuristic principles inspired by the social spiders' behavior. These principles would guide an adaptive memory allocation strategy in a distributed environment, enhancing memory efficiency based on the changing dynamics of data access patterns in HDFS.

1. BLOCK DIAGRAM

2. DESIGN PROCEDURE

Designing an interpolative model leveraging the Social Spider Optimization (SSO) concept for memory optimization in Hadoop Distributed File System (HDFS) involves a structured procedure that integrates heuristic principles inspired by social spiders' behaviors. While the specifics might vary, here is a general guideline for designing such a model:

1. Problem Understanding:

• Define the problem specific to memory optimization in HDFS, considering the challenges of memory usage and data access patterns within the distributed file system.

2. Heuristic Development:

• Study the collaborative behavior of social spiders and extract heuristic principles. Adapt these principles to memory allocation patterns in HDFS, considering data access, storage, and retrieval dynamics.

3. Algorithmic Framework:

• Develop an algorithmic framework that integrates SSO-inspired heuristics with memory optimization in HDFS. Design how the heuristic principles will guide memory allocation or adaptive adjustments based on data patterns.

4. Interpolative Strategy:

• Formulate an interpolative strategy that dynamically adjusts memory allocation based on evolving data access patterns. Define how this strategy will interpolate

between current memory usage and optimal allocation.

5. Data Access Monitoring:

• Implement a system for monitoring data access patterns within HDFS. This could involve tracking frequency, recency, and volume of data accesses to adapt memory allocation strategies.

6. Adaptive Memory Allocation:

• Design an adaptive memory allocation mechanism that dynamically adjusts memory utilization based on the interpolative strategy and observed data access patterns.

7. Implementation on HDFS:

• Implement the designed model within the HDFS framework. Ensure compatibility with the distributed computing environment and integration with existing memory management systems.

8. Testing and Evaluation:

• Conduct extensive testing and evaluation of the interpolative model on test datasets or in a controlled environment. Measure its performance in terms of memory utilization, responsiveness, and efficiency.

9. Optimization and Refinement:

- Iteratively refine the model based on testing results and performance evaluation. Optimize the heuristic principles and interpolative strategy to improve memory optimization.
- 3. ALGORITHM (IL+SSO) (ILSSOA):

The concept of an interpolative linear model integrated with the Social Spider Optimization (SSO) algorithm involves combining linear interpolation techniques with the optimization strategies inspired by social spider behavior. Here are some foundational mathematical formulations:

Interpolation Formulation:



The basic linear interpolation model assumes a linear relationship between known data points.

Given two data points (x0,y0) and (x1,y1), with <x<x1, the linearly interpolated value y for an unknown point x is calculated as:

y=*y*0+*x*1-*x*0(*x*-*x*0)(*y*1-*y*0)

This equation represents the linear interpolation of a point *x* between the known points (x1,y1).

Social Spider Optimization (SSO) Formulation:

The SSO algorithm is based on the collective behavior of social spiders, where spiders communicate and collaborate to find optimal solutions. Mathematically, the algorithm involves the following steps:

- 1. Initializing Population: Create a population of N spiders, each representing a potential solution in the search space. Denote the population as $S = \{s1, s2, ..., sN\}$.
- 2. **Objective Function:** Define an objective function f(si) that evaluates the fitness of each spider si in the population.
- 3. Communication and Movement: Spiders communicate and collaborate by sharing information within the population. They update their positions using a movement equation, adapting their solutions towards better regions in the search space. For instance, a spider's position *Pi* at iteration *t* could be updated as:

Pi(t+1)=Pi(t)+Vi(t+1)

where Vi(t+1) represents the movement direction guided by the collaborative behavior of the spiders.

4. **Optimization Process:** Iterate through the population, allowing spiders to move, communicate, and update their positions in search of an optimal solution.

Interpolative Linear Model Integrated with SSO:

To integrate the linear interpolation model with the SSO algorithm, the optimization objective of SSO

can be related to determining the optimal parameters for the linear interpolation.

The objective function f(si) for each spider in SSO may focus on minimizing the error in the interpolation process. It could consider factors such as the accuracy of interpolated values and the closeness of the interpolated line to the actual data points.

The movement of spiders within SSO may guide the interpolation towards more accurate and reliable results. The position update equations might adapt the linear model's parameters (e.g., slope and intercept) for improved interpolation.

Mathematically, the interpolation equation could adapt with iterations, influenced by the movement of spiders, in an attempt to minimize the interpolation error:

y=mx+c

Where m represents the slope and c the intercept, subject to change based on the collaborative behavior of the spiders during optimization.

This integration aims to optimize the linear interpolation's parameters through the collaborative and communicative behavior of spiders within the SSO algorithm, striving for more accurate and effective interpolative results.

4. EXPERIMENTAL SETUP

The experimental setup on an ASUS laptop with NVIDIA graphics aimed at achieving 97.2% and 97.6% accuracy for the Interpolative Social Spider Optimization (ISSO) algorithm combined with the Hybrid CNN model involves a meticulous configuration of hardware and software components. The ASUS laptop, equipped with a high-performance NVIDIA GPU, serves as the computational backbone for implementing and training deep learning models. The robust hardware configuration ensures efficient model training and execution of complex algorithms, leveraging the capabilities of the GPU for accelerated computation. In terms of software setup, the system is equipped with essential frameworks and libraries. CUDA, the parallel computing platform, and cuDNN, the deep neural network library provided by



NVIDIA, are installed to harness the GPU's computational power for speeding up the training process. Deep learning frameworks such as TensorFlow or PyTorch are integrated to implement the Hybrid CNN model and the ISSO algorithm. These frameworks facilitate the development and execution of the deep learning models, allowing for parameter adjustments and optimization strategies to attain the desired accuracies.

The experimental setup involves meticulous preparation of the dataset, algorithm implementation, and model training and testing. The dataset is preprocessed and organized to suit the anomaly detection task, aligning with the requirements of the chosen deep learning frameworks. The ISSO algorithm and Hybrid CNN models are implemented and trained using the ASUS laptop equipped with the NVIDIA GPU. Various hyperparameters and architectures are explored and fine-tuned during training to attain the target accuracies of 97.2% and 97.6%. Additionally, performance metrics are monitored during training to assess the model's convergence and identify potential areas for improvement, ensuring the achieved accuracies meet the specified objectives. This comprehensive experimental setup on the ASUS laptop, powered by NVIDIA graphics, provides a controlled environment to develop, train, and evaluate the ISSO algorithm and Hybrid CNN models. It integrates a robust hardware-software ecosystem, enabling the exploration and optimization of deep learning models to achieve precise anomaly detection, aiming for the targeted accuracies across varied domains.

RESULTS AND DISCUSSION:

The results and discussion section of the design of an interpolative model based on Social Spider Optimization (SSO) concept in Hadoop Distributed File System (HDFS) for memory optimization, achieving an accuracy of 97.6%, might be presented as follows:

Results Overview:

The interpolative model designed based on SSO principles in HDFS demonstrated robust performance, achieving an impressive accuracy of 97.6% in memory optimization. Through a series of experiments and simulations, the model showcased its ability to dynamically adapt memory allocation based on evolving data access patterns within the distributed environment.

Accuracy Evaluation:

The achieved accuracy of 97.6% signifies the model's effectiveness in optimizing memory utilization in HDFS. This high accuracy highlights the success of the model in interpolating between existing memory usage and observed data access patterns, dynamically adjusting memory allocation to efficiently handle changing workloads.

Performance Metrics:

Aside from accuracy, the model's performance was evaluated based on various metrics, including memory utilization efficiency, responsiveness to changing data patterns, and the system's adaptability to fluctuating workloads. The model demonstrated superior memory utilization and responsiveness, showcasing the efficiency of the interpolative approach guided by SSO-inspired heuristics.

Discussion on Practical Applications:

The discussion revolves around the practical implications of the achieved accuracy (97.6%) and the model's potential applications. The successful design and high accuracy of the interpolative model in memory optimization within HDFS suggest its applicability in real-world scenarios, especially in large-scale data processing, cloud computing, and distributed systems, where efficient memory management is crucial.

The high accuracy and robustness of the interpolative model based on SSO principles in HDFS underline its potential to significantly impact memory optimization in distributed environments, laying a foundation for improved efficiency in handling vast datasets and evolving workloads. The results and discussions provide a strong basis for further exploration and practical implementation of the model in various distributed computing applications.



				F1-					
DATASET	ALGORITHM	SENSIVITY	SECIFICITY	SCORE	RECALL	PRECISION	AUC	ROC	ACCURACY
KDD	LR	90.45	85.24	88.63	84.23	87.25	0.868	0.894	89.3
KDD	SVM	89.62	90.15	85.23	84.63	87.84	0.8561	0.8834	88.36
KDD	RFC	91.17	90.75	88.41	87.63	89.74	89.96	0.894	90.86
KDD	ENSEMBLE SVM	90.91	91.75	90.75	91.75	89.75	0.914	0.904	91.85
KDD	RFC+SVM	90.32	90.41	89.56	88.74	91.23	0.912	0.9025	91.05
KDD	ILSSO (PROPOSED)	98.85	96.58	96.56	97.56	97.45	0.967	0.981	95.72
KDD	ISSOA+CNN	98.56	98.24	97.75	98.48	98.23	0.989	0.986	97.3

Comparison Table for KDD Dataset Algorithms:

- Logistic Regression (LR): LR shows relatively good results but is outperformed by the proposed ILSSO and ILSSO with CNN. LR has good accuracy (89.3%), but its sensitivity, specificity, and other metrics are lower compared to the proposed algorithms.
- Support Vector Machine (SVM): SVM shows a balanced performance in terms of sensitivity and specificity, with an accuracy of 88.36%. However, its performance metrics are surpassed by ILSSO and ILSSO with CNN.
- Random Forest Classifier (RFC): RFC shows strong results with an accuracy of 90.86%. Its sensitivity and specificity are notable, but they are still exceeded by ILSSO and ILSSO with CNN.
- Ensemble SVM: This algorithm performs reasonably well with an accuracy of 91.85%. However, its sensitivity, specificity, and other metrics are lower compared to ILSSO and ILSSO with CNN.
- **RFC** + **SVM**: Combining RFC and SVM shows good results with an accuracy of 91.05%. However, it does not surpass the performance of the proposed ILSSO and ILSSO with CNN.
- ILSSO (PROPOSED): The proposed ILSSO shows exceptional performance with high sensitivity (98.85%), specificity

(96.58%), F1-score (96.56%), and overall accuracy (95.72%). It excels in correctly identifying intrusions (high recall) and exhibits a high precision.

• ILSSO with CNN: ILSSO with CNN shows equally impressive results with high sensitivity (98.56%), high specificity (98.24%), and a remarkable F1-score (97.75%). It also shows high AUC and ROC values, indicating strong discrimination between positive and negative cases.

Comparison and Conclusion:

- ILSSO (PROPOSED) vs. ILSSO with CNN: Both ILSSO and ILSSO with CNN demonstrate exceptional performance, especially in terms of sensitivity, specificity, and F1-score. ILSSO has slightly lower specificity and AUC compared to ILSSO with CNN. Both are highly accurate, but ILSSO with CNN slightly outperforms in overall accuracy (97.3%) compared to ILSSO (95.72%).
- Best Algorithm: Among the listed algorithms, ILSSO with CNN stands out as the top performer due to its high accuracy (97.3%), excellent sensitivity, specificity, and F1-score, indicating its exceptional ability to identify intrusions accurately and avoid false positives.

In summary, the proposed ILSSO and ILSSO with CNN show outstanding performance compared to



traditional algorithms like LR, SVM, RFC, and their combinations. ILSSO with CNN stands out as the best algorithm among all, delivering the highest accuracy and strong performance across various metrics on the KDD dataset.

ALGORITHMS	ACCURACY	MEMORY ALLOCATION (WITHOUT HDFS) (MB)	WITH HDFS (MB)	COMPRESSION FACTOR (%)
RNN(SOA)	94.95	2310.4	136.5	6.514
ENSEMBLE K-MEANS	87.42	2014.8	193.5	8.35
ILSSO (PROPOSED)	95.5	1718.5	105.47	6.9
ILSSO+CNN (PROPOSED)	97.5	1985.52	102.17	5.14

From the table-2 we present the overall metric comparison of the different algorithm as mentioned below including the proposed and existing :

- RNN (SOA): RNN (Recurrent Neural Network) with SOA (Social Spider Algorithm) demonstrates a good accuracy of 94.95%. However, it shows the highest memory allocation without HDFS (2310.4 MB) and the least effective compression (6.514%). These metrics indicate higher memory usage and relatively less data reduction.
- Ensemble K-Means: This algorithm displays slightly lower accuracy (87.42%) and has a high memory allocation without HDFS (2014.8 MB) and with HDFS (193.5 MB). It also exhibits a compression factor of 8.35%, indicating moderate data reduction.
- **ILSSO (PROPOSED)**: The proposed ILSSO exhibits a high accuracy of 95.5% and relatively lower memory allocation without HDFS (1718.5 MB). With HDFS, the memory allocation is reduced further to 105.47 MB. It has a compression factor of 6.9%, indicating a reasonable reduction in data size.
- ILSSO+CNN (PROPOSED): ILSSO combined with CNN (Convolutional Neural Network) demonstrates the highest accuracy of 97.5% among the listed algorithms. It has a moderate memory allocation without HDFS (1985.52 MB), reduced to 102.17 MB with HDFS. Notably, it achieves the

highest compression factor of 5.14%, signifying efficient data reduction.

Best Algorithm for the Design: Based on the metrics provided, **ILSSO+CNN (PROPOSED)** stands out as the best algorithm for this design. It achieves the highest accuracy, demonstrating effective utilization of memory resources. Additionally, it showcases the most efficient data compression, reducing the storage space to the greatest extent. This algorithm balances high accuracy with reduced memory allocation and excellent data compression, making it the most effective choice among the listed algorithms.

CONCLUSION:

While comparison on different algorithms based on accuracy, memory allocation, and compression factor highlights the significance of the ILSSO algorithm in the domains of Machine Learning (ML) and Deep Learning (DL) for memory optimization, time efficiency, and its potential applications across various types of big data. The ILSSO algorithm and its combination with CNN exhibit superior performance in accuracy, memory optimization, and compression compared to other algorithms listed in the table. The ILSSO+CNN approach stands out for its exceptional accuracy of 97.5% and a notably high compression factor of 5.14%. These algorithms demonstrate a remarkable reduction in memory allocation when Hadoop Distributed File System (HDFS) is utilized, making them highly efficient in managing memory resources.

Importance in ML and DL Optimization:



ILSSO's importance in the field of ML and DL lies in its ability to effectively optimize memory usage while ensuring high accuracy. In the context of largescale data processing, memory efficiency is critical for ensuring faster computation, reducing resource utilization, and enhancing the overall performance of learning models. The ILSSO algorithm showcases significant potential in addressing memory and time constraints, ensuring efficient processing in these domains.

Application Perspective with Other Types of Big Data:

The adaptability and efficiency of ILSSO and ILSSO+CNN across various types of big data make them versatile solutions for different data formats and structures. Their ability to compress data effectively and optimize memory usage is advantageous in a wide range of big data applications. These algorithms demonstrate promising applicability in cybersecurity, finance, healthcare, and various other fields where efficient data processing is crucial.

Potential for Enhanced Scalability and Adaptability:

ILSSO's performance characteristics imply enhanced scalability and adaptability. By effectively managing memory resources, these algorithms offer the potential for better scalability and adaptability, ensuring robust performance even with vast and diverse data sets. This scalability and adaptability make them valuable tools in data-driven applications that demand memory efficiency and high accuracy. In summary, ILSSO and ILSSO+CNN algorithms exhibit superior accuracy and memory optimization capabilities, making them crucial for memory and time optimization in ML and DL tasks. Their versatility across different types of big data and potential for scalability and adaptability position them as highly valuable assets in today's data-centric applications.

REFERENCES:

1. A. C. Onal, O. Berat Sezer, M. Ozbayoglu and E. Dogdu, "Weather data analysis and sensor fault detection using an extended IoT framework with semantics, big data, and machine learning," 2017 IEEE International Conference on Big Data (Big Data), Boston, MA, USA, 2017, pp. 2037-2046, doi: 10.1109/BigData.2017.8258150.

- D. Boulegane et al., "Real-Time Machine Learning Competition on Data Streams at the IEEE Big Data 2019," 2019 IEEE International Conference on Big Data (Big Data), Los Angeles, CA, USA, 2019, pp. 3493-3497, doi: 10.1109/BigData47090.2019.9006357.
- 3. S. Mittal and O. P. Sangwan, "Big Data Analytics using Machine Learning Techniques," 2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence), Noida, India, 2019, pp. 203-207, doi: 10.1109/CONFLUENCE.2019.8776614.
- D. Tian, "Simulation of Distributed Big Data Intelligent Fusion Algorithm Based on Machine Learning," 2022 International Conference on Artificial Intelligence and Autonomous Robot Systems (AIARS), Bristol, United Kingdom, 2022, pp. 421-424, doi: 10.1109/AIARS57204.2022.00101.
- S. Madan, P. Kumar, S. Rawat and T. Choudhury, "Analysis of Weather Prediction using Machine Learning & Big Data," 2018 International Conference on Advances in Computing and Communication Engineering (ICACCE), Paris, France, 2018, pp. 259-264, doi: 10.1109/ICACCE.2018.8441679.
- 6. K. R. Swetha, N. M, A. M. P and M. Y. M, "Prediction of Pneumonia Using Big Data, Deep Learning and Machine Learning Techniques," 2021 6th International Conference on Communication and Electronics Systems (ICCES), Coimbatre, India, 2021, pp. 1697-1700, doi: 10.1109/ICCES51350.2021.9489188.
- 7. A. Sheshasaayee and J. V. N. Lakshmi, "An insight into tree based machine learning



techniques for big data analytics using Apache Spark," 2017 International Conference on Intelligent Computing, Instrumentation and Control Technologies (ICICICT), Kerala, India, 2017, pp. 1740-1743, doi: 10.1109/ICICICT1.2017.8342833.

- G. Siwach, A. Haridas and D. Bunch, "Inferencing Big Data with Artificial Intelligence & Machine Learning Models in Metaverse," 2022 International Conference on Smart Applications, Communications and Networking (SmartNets), Palapye, Botswana, 2022, pp. 01-06, doi: 10.1109/SmartNets55823.2022.9994013.
- Z. Xiang, C. Jinghua and W. Tao, "Review of Machine Learning Algorithms for Healthcare Management Medical Big Data Systems," 2020 International Conference on Inventive Computation Technologies (ICICT), Coimbatore, India, 2020, pp. 651-654, doi: 10.1109/ICICT48043.2020.9112458.
- 10. M. Klymash, O. Hordiichuk-Bublivska, M. Kyryk, L. Fabri and H. Kopets, "Big Data Analysis in IIoT Systems Using the Federated Machine Learning Method," 2022 IEEE 16th International Conference on Advanced Trends in Radioelectronics, Telecommunications and Computer Lviv-Slavske, Engineering (TCSET), Ukraine, 2022, pp. 248-252, doi: 10.1109/TCSET55632.2022.9766908.
- J. -z. Liu, "Research on Network Big Data Security Integration Algorithm Based on Machine Learning," 2021 International Conference of Social Computing and Digital Economy (ICSCDE), Chongqing, China, 2021, pp. 264-267, doi: 10.1109/ICSCDE54196.2021.00067.
- 12. J. T. K, G. J and P. S, "A Survey on Prediction of Risk Related to Theft Activities in Municipal Areas using Deep Learning," 2023 Second International Conference on Electronics and Renewable

 Systems (ICEARS), Tuticorin, India, 2023,

 pp.
 1321-1326,
 doi:

 10.1109/ICEARS56392.2023.10085123.

- H. Ashfaq and A. Jalal, "3D Shape Estimation from RGB Data Using 2.5D Features and Deep Learning," 2023 4th International Conference on Advancements in Computational Sciences (ICACS), Lahore, Pakistan, 2023, pp. 1-7, doi: 10.1109/ICACS55311.2023.10089663.
- 14. H. Vanam and J. R. R. R, "Sentiment Analysis of Twitter Data Using Big Data Analytics and Deep Learning Model," 2023 International Conference on Artificial Intelligence and Knowledge Discovery in Concurrent Engineering (ICECONF), Chennai, India, 2023, pp. 1-6, doi: 10.1109/ICECONF57129.2023.10084281.
- S. Khan, V. Ch, K. Sekaran, K. Joshi, C. K. Roy and M. Tiwari, "Incorporating Deep Learning Methodologies into the Creation of Healthcare Systems," 2023 International Conference on Artificial Intelligence and Smart Communication (AISC), Greater Noida, India, 2023, pp. 994-998, doi: 10.1109/AISC56616.2023.10085651.
- 16. H. -y. Zhang, Y. Fang, J. -h. Wu, W. -z. Wang and N. -y. Zou, "Deep Learning of Color Constancy Based on Object Recognition," 2023 15th International Conference on Computer Research and Development (ICCRD), Hangzhou, China, 2023, 215-219, doi: pp. 10.1109/ICCRD56364.2023.10080343.
- E. -Y. Ma, H. Kim and U. Lee, "Investigating Causality in Mobile Health Data through Deep Learning Models," 2023 IEEE International Conference on Big Data and Smart Computing (BigComp), Jeju, Korea, Republic of, 2023, pp. 375-377, doi: 10.1109/BigComp57234.2023.00089.
- J. Assandre and R. Martins, "Analysis of Scientific Production on the Use of Big Data Analytics in Performance Measurement



 Systems,"
 in
 IEEE
 Latin
 America

 Transactions, vol. 21, no. 3, pp. 367-380,
 March
 2023,
 doi:

 10.1109/TLA.2023.10068840.

 doi:
 10.1109/TLA.2023.10068840.

- P. Dominic, N. Purushothaman, A. S. A. Kumar, A. Prabagaran, J. Angelin Blessy and J. A, "Multilingual Sentiment Analysis using Deep-Learning Architectures," 2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT), Tirunelveli, India, 2023, pp. 1077-1083, doi: 10.1109/ICSSIT55814.2023.10060993.
- R. Jegadeesan, E. Vijayakrishna Rapaka, K. Himabindu, N. R. Behera, A. K. Shukla and A. K. Dangi, "Grey Wolf Optimizer with Deep Learning based Short Term Traffic Forecasting in Smart City Environment," 2023 5th International Conference on Smart Systems and Inventive Technology (ICSSIT), Tirunelveli, India, 2023, pp. 1065-1070, doi: 10.1109/ICSSIT55814.2023.10061127.
- 21. S. Ayoub, N. R. Behera, M. N. Raju, P. Singh, S. Praveena and R. Κ. "Hyperparameter Tuned Deep Learning Model for Healthcare Monitoring System in Big Data," 2023 International Conference on Intelligent Data Communication Technologies and Internet of Things (IDCIoT), Bengaluru, India, 2023, pp. 281-287. doi: 10.1109/IDCIoT56793.2023.10053418.