

conVergence

INTERNATIONAL JOURNAL OF ICTAK
A MULTIDISCIPLINARY JOURNAL OF ENGINEERING, TECHNOLOGY & EMPLOYABILITY



A Govt. of India supported, Govt. of Kerala partnered Social Enterprise

VOLUME 07

ISSUE 01

DEC 2022



Culture eats strategy for breakfast...

~ Peter Drucker



Editorial Board

Chief Patron

Dr. P.V. Unnikrishnan

Member Secretary

Kerala Development and Innovation Strategy Council (K-DISC)

Chief Editor

Dr. Manoj A.S.

Head - Planning, Competency Development & Innovations

ICT Academy of Kerala (ICTAK)

Editorial Advisory Board

Dr. Saji Gopinath

Vice-Chancellor, Digital University Kerala

Prof. Mithileswar Jha

Retd. Professor of Marketing, IIM B

Dr. A. Damodaran

Retd. Professor, IIM B

Distinguished Professor, ICRIER, New Delhi

Prof. L.S. Murthy

Professor, IIM B

Prof. Prahalad Vadakkepat

Associate Professor, NUS

Dr. Rajasree M.S.

Joint Director, Department of Technical Education



Contents

1. COVID-19 Future Forecasting Using Bagging Method and Outbreak Using Multivariate Regression	7-14
2. Edge Computing: Technology Born For Elegant Connectivity	15-21
3. Online English Teaching Using Artificial Intelligence Under Big Data Environment	23-31
4. Employee Attrition Prediction Using Various Machine Learning Algorithms	33-40
5. An Authenticatable (2, 3) Secret Sharing Scheme Using Meaningful Share Images Based on Hybrid Fractal Matrix	41-54
6. Physics Equation of Motion Problem Solver	55-61
7. Pneumonia Category Detection Using Deep Learning	63-69



Introduction

Dear Readers,

I am delighted to present to you the 7th issue of the ICT Academy of Kerala (ICTAK) International Journal 'Convergence'. This edition holds a special significance as it explores the intersection of 'Skills, Engineering & Employability' with a focus on the emerging theme of the Metaverse. The Metaverse has emerged as a transformative concept that encompasses virtual reality, augmented reality, and other immersive technologies. It has the potential to revolutionize how we live, work, and interact, opening up new opportunities and challenges for individuals, businesses, and society as a whole. This issue of 'Convergence' aims to delve into the implications of the Metaverse on skills, engineering education, and employability.

Our esteemed contributors, comprising researchers, academicians, industry experts, and visionaries, have shared their insights and research findings through scholarly articles, case studies, and thought-provoking editorials. Their collective knowledge and expertise provide a comprehensive exploration of the Metaverse's impact on skills development, engineering education, and the evolving nature of employability. The articles featured in this issue cover a wide range of topics, including, The Metaverse and Skill Development, Engineering Education in the Metaverse, Workforce Transformation, Ethical Considerations & Entrepreneurship and Innovation. By exploring the convergence of Skills, Engineering & Employability with the Metaverse, this issue aims to contribute to the ongoing discourse on preparing individuals and organizations for the future of work and technological advancements.

I would like to express my deepest gratitude to all the authors who have contributed their valuable articles to this edition. Their expertise, research, and insights have enriched the content and provided readers with a diverse range of perspectives. I would also like to extend my appreciation to the editorial team, reviewers, and the entire publishing committee for their tireless efforts in ensuring the quality and relevance of the articles included in this edition. Lastly, I extend my sincere thanks to our readers for their continued support and engagement. It is your enthusiasm and intellectual curiosity that motivates us to bring forth thought-provoking content and foster a thriving academic and professional community. I invite you to immerse yourself in the pages of this 7th issue of 'Convergence' and explore the enlightening articles on Skills, Engineering & Employability in the context of the Metaverse. May these articles inspire you, spark discussions, and propel us towards a future that harnesses the full potential of the Metaverse for the betterment of society. Thank you for your unwavering support.

Dr. Manoj A.S.

Chief Editor, Convergence

COVID-19 Future Forecasting Using Bagging Method and Outbreak Using Multivariate Regression

Sona V. Jose, Francy T.L.

Department of Computer Science Vimala College Thrissur, Kerala |sonavjose2000@gmail.com

Department of Computer Science Vimala College Thrissur, Kerala |francyjai@gmail.com

Abstract

The COVID-19 pandemic has led to an immense loss of human life all around the world. Our social life and economic worth have been severely damaged due to this pandemic. During the covid-19 pandemic in various countries, many authorities have used machine learning techniques to make accurate forecasting for the covid-19 pandemic. Machine learning (ML) is a type of artificial intelligence (AI) that allows software programs to improve prediction accuracy. This study demonstrates the capability of different ML models to forecast the number of upcoming patients affected by COVID-19 worldwide. The most affected countries are the US, Spain, Italy, France, Germany, Russia, Iran, the United Kingdom, Turkey, and India. Our work focus on the bagging method for future forecasting and multivariate regression has been used to predict the threatening factors of covid-19 and we attain an accuracy of 99.99%. Experiments have proved among these prediction models, the improved bagging method and multivariate regression can predict the future development trend of the global pandemic more precisely.

Keywords: COVID-19, Machine learning, Bagging, Multivariate, Forecast

1. Introduction

The World Health Organization (WHO) was informed of cases of pneumonia of unknown cause in Wuhan City, China On 31 December 2019, The new coronavirus was identified by the Chinese government on January 7, 2020, as the cause and was temporarily named "2019nCoV". In a short period of time, COVID-19 began to spread rapidly, affecting almost the entire world. Almost 216 countries around the world had fought the coronavirus pandemic.

COVID-19 virus will experience mild to moderate respiratory illness and recovery without requiring special treatment in most of the infected people. To prevent the infection socially, a lockdown and a travel ban were imposed in almost every country, which resulted in the halt of all financial and social activity in society. This drive to diminish global supply chains badly results in the global economy being in the worst condition. In India, the government also imposed a nationwide lockdown for the first time in history from March 22, 2020,

and continued it up to May 30, 2020. With the spread of Covid-19, the government has restricted all sectors including, education, manufacturing, hospitality, services, and transportation. During lockdown, people are forced to start working from home. School and college classes are running online, and an evident number of people shifted to a digital platform. On May 30, 2020, the outbreak of this coronavirus disease worsened, including 5,704,736 confirmed cases and 357,736 confirmed deaths in 216 countries. There are also 165,799 confirmed cases in India, with 4,706 casualties.

Older people and people who have medical problems like chronic respiratory disease, cardiovascular disease, diabetes, and cancer are more vulnerable. The best way to prevent and reduce infection is to pay close attention to the COVID-19 virus, the disease it causes, and how it spreads. Protect yourself and others from infection by washing your hands frequently or using disinfectants to minimize social gatherings. It is now very clear that COVID-19 has a very devastating impact around the world. People are panicked, emotionally anxious, and depressed, and relied on news sources for COVID-19 symptoms and their prevention by treatment.

According to reports from various countries, the number of cases of infected people has increased dramatically. This study helps the general public and government to plan and determine strategies for combating the coronavirus. At the early stage of the outbreak, we have only a limited amount of resources to predict the impact of coronavirus in the upcoming days. Now, much information is available for the prediction because the virus has spread on a large scale. Every day thousands of people were affected by the coronavirus in most countries and the death rate also increases accordingly. The coronavirus spread through physical contact, respiratory droplets, being near infected persons, and touching contaminated surfaces. The main reason for the outbreak of the Covid-19 virus is that infected people have been

asymptomatic for days.

However, different models have different predictions, making it difficult to make the right choices when performing the right actions. We combined the relevant experiences of SARS in the past, the infectious disease dynamic model, and the existing data on COVID-19 to compare the forecast of future epidemic situations. The paper has been done to determine the most suitable prediction model for this epidemic. The bagging method is used to predict the upcoming 30 days of confirmed cases, death cases, and recovery cases. I used multivariate regression to predict the end date of this pandemic. The learning models have been trained using the COVID-19 patient stats dataset provided by Johns Hopkins University. For the analysis, we have divided the datasets of confirmed, recovered, and death cases into training and testing. The results show that Bagging and Multivariate regression produced better forecasting compared to other methods.

2. Related Work

After a thorough investigation and investigation, we found a similar type of study, but more noteworthy. "COVID-19 future forecasting the usage of supervised device learning fashions" [1] a magazine written by honourable professors, has used comparable techniques and fashions, however with greater studies and experimentation. we are able to now not be evaluating our paintings with theirs as they have used fewer facts compared to my work.

Saudi Sheikh et al. [2] Presented a paper explaining the use of various retrospective models that predicts the spread of covid-19 in India. Linear and Polynomial models have different degrees and are designed to predict the test data. Also found that the polynomial model with 5 degrees is usually better aligned with real values. We found that the higher the ordinal value, the higher the accuracy and the minimum mean square error. Therefore, the polynomial regression model has a higher performance and a line model than other

degree models. Performance tests were performed using Mean Total Performance Error (MAPE) and R-squared (R²). They used heatmaps to identify correlations between organizations. To get data prediction, they do Tableau use. The results obtained are satisfactory and can be overcome by increasing the size of the data and using better algorithms.

Ajinkya Kunjir et al. [3] in their paper analyzed the WHO Database predicting the spread of COVID-19. Use algorithms such as Short-Term Memory (LSTM), Convolutional Neural Network (CNN), and Decision Tree. The database contained 91 data days. The dataset was over and only India, China, and Canada were studied. Each algorithm was present and tested in each country and compared actual features (confirmed cases, death, recovery). The test was performed for an unknown 25 time series days. It has been observed that the CNN model fits well and is familiar with overrides the actual database curve. Based on R² points, CNN did very well followed by Decision Tree and LSTM finally, in all three countries. The model performed well and provided accurate results. Many countries and regions and a few other algorithms can be added to better read by comparison. The models can be tested for a few additional function metrics like MSE, and Variance.

Saksham Gera et al. [4] did some research by predicting Covid19 styles using various machine learning algorithms. The database was used in this study for 292 days. Many algorithms are the same Down the line, the nearest K neighbour, Vector support Machine, Random Forest, and Elastic Smoothing were used lesson. The paper also evaluates the performance of all algorithms implemented using the Root Mean Square Error (RMSE), R School Square (R² School), and Total Error Means (MAE). The result was analyzed in two different respects. 1) Predictability of newly infected cases 2) Death prediction charges. To predict new cases and ES deaths performance was better than all other algorithms and SVM has malfunctioned. The reason for SVM the poor performance was due to the presence of outliers inside

data and fluctuations in data set values. My performance would have been better if the outliers would have been properly treated during the data purification process.

H Zakiyyah and S Suyanto [5] in their paper have developed Machine learning models to predict the newly infected cases of Covid 19 and deaths in Indonesia. The authors have trained their models on the Covid-19 Indonesia Time Series All Dataset (CITSAD) of ten provinces. Three prediction models were made using Gaussian Naïve Bayes (GNB), Support Vector Machine (SVM), and Decision Tree (DT) algorithms. The performance evaluation was done based on accuracy and processing time. Based on R² score performance, the Decision Tree (DT) model proved to be the best model for the prediction as it had the highest accuracy and also the least processing time out of all the models being tested. The reason behind the success of the Decision tree (DT) was that the model could split the complex decision-making process into simpler ones making it more effective than other algorithms for more accurate prediction.

3. Materials and methodology

3.1. Data set

This segment describes the data set we've got used to expecting the covid 19 instances of followed countries, and globally focusing on the wide variety of latest fantastic instances, the wide variety of deaths, and the wide variety of recoveries. The Centre for Systems Science and Engineering, Johns Hopkins University offers the data set used in this study. The repository became broadly speaking made to be had for the visible dashboard of 2019 Novel Corona virus via way of means of the college and became supported via way of means of the ESRI Living Atlas Team. The folder consists of day-by-day time collection precis tables, consisting of the wide variety of shown instances, deaths, and recoveries. All statistics are from the day-by-day case record and the replace frequency of statistics is one day.

3.2. Methodology

The proposed methodology is a two-step process that aims

A. Prediction and future forecasting of covid cases including death, recovery, and confirmed cases by using bagging models.

a. Bagging Method

A Bagging classifier is an ensemble meta-estimator that fits base classifiers on arbitrary subsets of the authentic data set after which aggregate their particular prognostications (either by way of voting or via averaging) to form a final prediction. This type of meta-estimator can typically be used as a way to lessen the friction of a black-field estimator (e.g., a decision tree), by introducing randomization into its creation fashion after which making an ensemble out of it. Every base classifier is trained in resemblance with a training set that's generated by aimlessly drawing, with a cover, N exemplifications (or statistics) from the authentic training data set – in which N is the size of the unique training set. The training set for each of the base classifiers is unbiased from every other. Most of the authentic records may be repeated in the preceding training set whilst others may be left out.

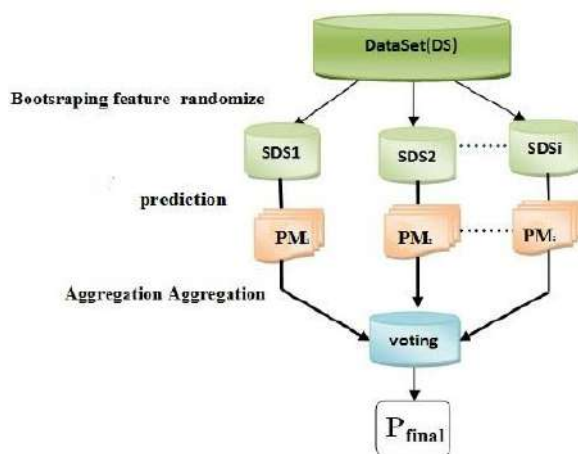


Fig. 1: Bagging (Bootstrap Aggregation)

b. Algorithm for the Bagging classifier

i) Classifier generation:

Let N be the size of the training set for each of t iterations:

Sample N instances with replacement from the original training set. Apply the learning algorithm to the sample.

Store the resulting classifier.

ii) Classification:

For each of the t classifiers:

Predict the class of instance using the classifier.

Return class that was predicted most often.

B. Covid-19 outbreak prediction with Multivariate regression techniques.

a) Multivariate linear regression

Regression analysis is one of the most popular ways to evaluate facts. It is based on the gadget mastering algorithm that is being tracked. Regression analysis is a basic statistical approach that allows you to examine relationships or additional variables between variables in a dataset. Simple linear regression is a regression model that estimates the relationship between an established variable and an independent variable using a straight line. However, multiple regression estimates the relationship between two or more unbiased variables and one dependent variable. The difference between these two modes is a fair range of variables. Multivariate regression is a machine that is constantly monitored and evaluates a set of rules that reference several fact variables. Multivariate regression is an extension of multiple regression with one structured variable and multiple unbiased variables. Primarily expect output based on various independent variables. Multivariate regression aims to create an expression that can explain how a variable

responds to changes in other variables at the same time.

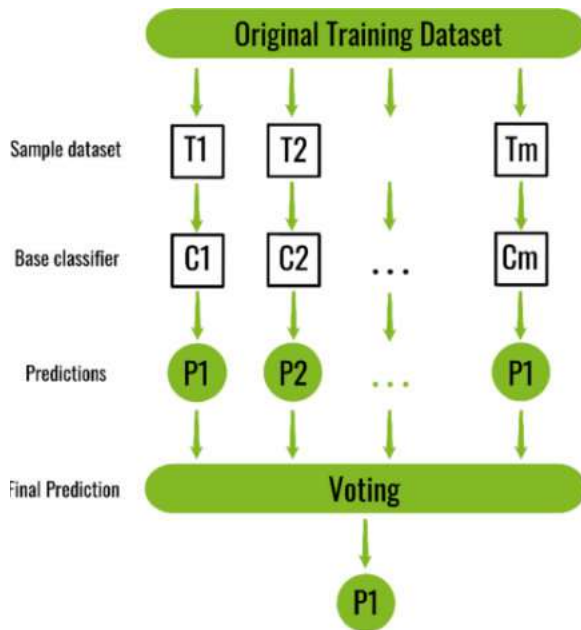


Fig. 2: Bagging Classifier

The simple regression linear model represents a straight line meaning y is a function of x . When we have an extra dimension (z), the straight line becomes a plane. Here, the plane is the function that expresses y as a function of x and z . The linear regression equation can now be expressed as:

$$y = m_1.x + m_2.z + c$$

y is the dependent variable, that is, the variable that needs to be predicted. x is the first independent variable. It is the first input. m_1 is the slope of x . It lets us know the angle of the line (x). z is the second independent variable. It is the second input. m_2 is the slope of z . It helps us to know the angle of the line (z). c is the intercept. A constant that finds the value of y when x and z are 0. The equation for a model with two input variables can be written as:

$$y = \beta_0 + \beta_1.x_1 + \beta_2.x_2$$

What if there are three variables as inputs? Human visualizations can be only three dimensions. In the machine learning world, there can be n number of dimensions. The equation for a model with three input variables can be written as:

$$y = \beta_0 + \beta_1.x_1 + \beta_2.x_2 + \beta_3.x_3$$

Below is the generalized equation for the multivariate regression model-

$y = \beta_0 + \beta_1.x_1 + \beta_2.x_2 + \dots + \beta_n.x_n$ Where n represents the number of independent variables, β_0 β_n represents the coefficients, and x_1 x_n is the independent variable. The multivariate model helps us in understanding and comparing coefficients across the output. Here, the small cost function makes Multivariate linear regression a better model.

b) Steps of multivariate regression evaluation

Steps concerned in Multivariate regression evaluation are characteristic choice and characteristic engineering, normalizing the capabilities, deciding on the loss feature and speculation, placing speculation parameters, minimizing the loss feature, checking out the speculation, and producing the regression version.

Feature choice- The choice of capabilities is a vital step in multivariate regression. Feature selection is also known as variable selection. It will become vital for us to select sizable variables for higher version building.

Normalizing Features- We want to scale the capabilities because it continues widespread distribution and ratios in data. This will result in a green evaluation. The cost of every characteristic also can be changed.

Select Loss feature and Hypothesis- The loss feature predicts every time there may be an error. In other words, speculation prediction differs from actual values. The projected cost from the characteristic/variable is the subject of speculation here.

Set Hypothesis Parameters- The speculation parameter wishes to be set in the sort of manner that reduces the loss feature and predicts well.

Minimize the Loss Function- The loss feature wishes to be minimized with the aid of the use of a loss minimization set of rules at the dataset if you want to assist in adjusting speculation parameters. After the loss is minimized, it may be used for similar actions. Gradient descent is one of the algorithms generally used for loss minimization.

Test the speculation feature- The speculation feature wishes to be checked as well, as its miles predicting values. Once that is done, it must be examined to take a look at the data.

4. Results and discussion

4.1. Forecasting framework

This effort seeks to build a system for future forecasting of the number of patients affected by COVID-19 using machine learning approaches. The dataset used for the study includes daily statistics on the number of newly infected patients, deaths, and recoveries caused by Covid-19 around the world. The world is undergoing

an alarming scenario in which the death rate and confirmed cases are increasing day by day. This paper helps to forecast the number of people that can be affected in terms of deaths and new confirmed cases including the number of expected recoveries for the upcoming days. It also predicts that it is a promising method for the current scenario of the COVID-19 pandemic, as well as the prediction of the pandemic deadline. We used bagging and machine learning model multivariate regression to predict the number of newly infected cases, the number of recoveries, the number of deaths, and the deadline.

4.2. Forecasting of recovered cases

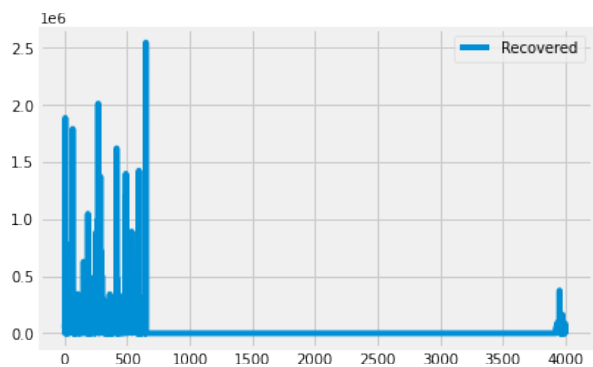


Fig. 4: Recovered cases

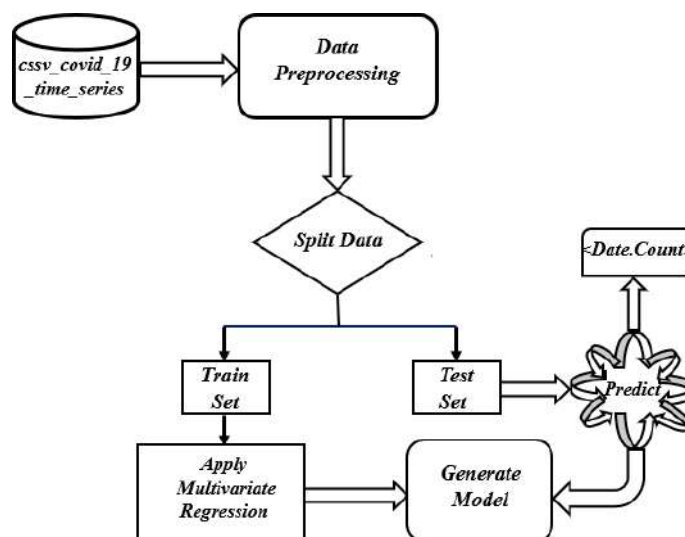


Fig. 3: Framework

To forecast and analyze the recovery rate of disease, we've finished an evaluation and have a look at the following fashions and the usage of recuperation records of the ten followed countries, and worldwide. We can see from the trends that recovery data is also non-stationary. So, we have performed stationarity techniques similar to as discussed in the section to evaluate the bagging method.

```
accuracy :
99.99999970892587
Recovery cases count as per bagging ensembles:
34355509
```

Fig. 5: Forecasting of recovered cases

4.3. Forecasting of Death Cases

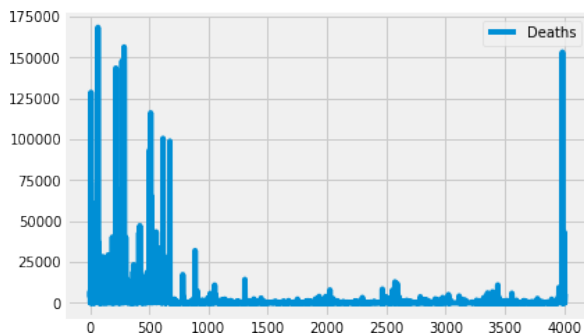


Fig. 6: Death cases

Corona virus has taken many lives. So, it is necessary to analyze the fatality rate of the virus, and forecasting to highlight future cases which can guide governments to act in advance. In this section, we have evaluated the forecasting models for death cases in the adopted countries, and worldwide. We have converted the non-stationary fatality data into stationary form

```
accuracy :
99.99997829818724
Death cases count as per bagging ensembles:
460791
```

Fig. 7: Forecasting of death cases

to fit the bagging model similar to as discussed in the above section.

4.4. Forecasting of new infected confirmed cases

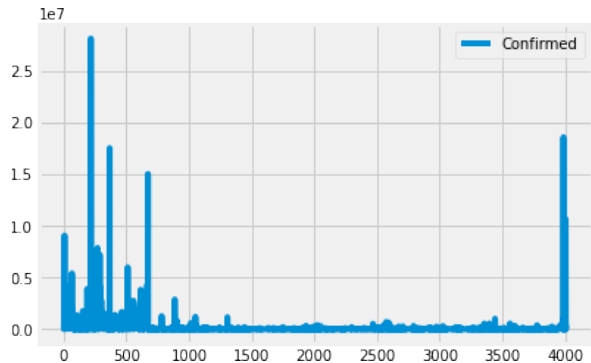


Fig. 8: Confirmed case

In this section, we have highlighted the fitted accuracy of the model using confirmed cases. We have evaluated the forecasting models for confirmed cases in the adopted countries, and worldwide. We have converted the non-stationary fatality data into stationary form to fit the bagging model similar to as discussed in the above section

```
accuracy :
99.99999941785174
confirmed cases count as per bagging ensembles:
34355509
```

Fig. 9: Forecasting of confirmed cases

4.5. Outbreak prediction

```
Approximate Days to Outbreak 728
<class 'int'>
Outbreak date 2024-03-29 13:33:00
```

Fig. 10: Outbreak prediction

To covid 19 pandemic in various countries, many authorities have used machine learning techniques to make accurate

forecasting for the covid 19 pandemic. The improved Multivariate regression can predict the future development trend of the global pandemic more precisely.

5. Conclusion

The work proposed here is intended to make accurate forecasting for the covid 19 pandemic. This proposed model is based on machine algorithms and the Ensemble method that have been tested on the COVID-19 dataset of The Centre for Systems Science and Engineering, Johns Hopkins University. While other researchers have tried machine learning methods before, they tend to achieve only decent accuracy of around 70-80%, but COVID-19 is a case which is crucial to detect the number of upcoming cases. Here, False reports are more dangerous, and so machine learning applications must deliver more accuracy in order to avoid false decisions and provide true efficacy as predicting tools. When we compared our proposed work with existing schemes, it

became evident that the 99.99 % accuracy we achieved was significantly better than most other machine learning algorithms and models have managed. Various evaluations of our proposed model, helped us to select the most suitable machine learning algorithm and ensemble model to accomplish the desired task.

6. Future work

In future work, we will use other Ensemble techniques with the proposed work to increase both its efficiency and accuracy.

7. Acknowledgement

We are extremely grateful to God Almighty, whose blessings have given us the courage and strength to complete this work successfully. We would like to express our sincere gratitude to all the teaching staff, for their valuable guidance and support at each stage of the work. We express our heartiest gratitude to my parents for the support given in connection with the work.

8. References

F. Rustam, A. A. Reshi, A. Mehmood, S. Ullah, B.-W. On, W. Aslam, and G. S. Choi, "Covid-19 future forecasting using supervised machine learning models," *IEEE access*, vol. 8, pp. 101489–101499, 2020.

M. Singh and S. Dalmia, "Prediction of number of fatalities due to covid-19 using machine learning," in *2020 IEEE 17th India Council International Conference (INDICON)*, pp. 1–6, IEEE, 2020.

V. K. Gupta, A. Gupta, D. Kumar, and A. Sardana, "Prediction of covid-19 confirmed, death, and cured cases in India using random forest model," *Big Data Mining and Analytics*, vol. 4, no. 2, pp. 116–123, 2021.

J. A. Webb and J. K. Aggarwal, "Visually interpreting the motion of objects in space," *Computer*, vol. 14, no. 08, pp. 40–46, 1981.

H. Zakiyyah and S. Suyanto, "Prediction of covid-19 infection in Indonesia using machine learning methods," in *Journal of Physics: Conference Series*, vol. 1844, p. 012002, IOP Publishing, 2021.

I. Arpacı, S. Huang, M. Al-Emran, M. N. Al-Kabi, and M. Peng, "Predicting the covid-19 infection with fourteen clinical features using machine learning classification algorithms," *Multimedia Tools and Applications*, vol. 80, no. 8, pp. 11943–11957, 2021.

D. Assaf, Y. Gutman, Y. Neuman, G. Segal, S. Amit, S. Gefen-Halevi, N. Shilo, A. Epstein, R. Mor-Cohen, A. Biber, et al., "Utilization of machine-learning models to accurately predict the risk for critical covid- 19," *Internal and emergency medicine*, vol. 15, no. 8, pp. 1435–1443, 2020.

K. B. Prakash, S. S. Imambi, M. Ismail, T. P. Kumar, and Y. Pawan, "Analysis, prediction and evaluation of covid-19 data sets using machine learning algorithms," *International Journal*, vol. 8, no. 5, pp. 2199–2204, 2020.

A. F. de Moraes Batista, J. L. Miraglia, T. H. R. Donato, and A. D. P. Chiavegatto Filho, "Covid-19 diagnosis prediction in emergency care patients: a machine learning approach," *MedRxiv*, 2020.

Y. Zoabi and N. Shomron, "Covid-19 diagnosis prediction by symptoms of tested individuals: a machine learning approach," *MedRxiv*, 2020.



Edge Computing: Technology Born For Elegant Connectivity

Monte Pious

Assistant Professor, Department of Computer Applications, Michael's Institute of Management and Technology, Mayithara | E-mail:montepious@gmail.com

Abstract

Edge computing technology brings the applications close to the user while keeping some portions at the server side. Edge computing architecture consists of three layers. The implementation of Edge Computing comprises five phases. Edge computing has imperative advantages like demanding only low latency. Edge computing also faces some challenges like various reliability and security concerns. Edge computing has profuse applications. Edge computing is very effective for Internet of Things (IoT) and Fifth Generation Mobile Network (5G). Multi Access Edge Computing and Fog Computing are some of the prominent trends in the Edge Computing field.

Keywords: *Edge Computing, Cloudlet, Edge Data centers, Multi Access Edge Computing, Fog Computing*

1. Introduction

As day by day the amount of data needed to handle and process is increasing expeditiously. Also, IoT devices are becoming emerging and prevalent. There comes the scope of Edge Computing technology. Edge computing is a proximity solution which delivers the data and essential computing needed from the cloud or any remote data center to the edge of the network in order to bring the services closer to the end user. It reduces the processing workload on backend servers and network latency substantially.

2. Edge Computing and Architecture

Edge computing is an emerging computing paradigm which optimizes IOT devices and applications where user data is processed

at the periphery of the network. Edge Computing brings computing in close proximity to the origin of the data.

Let us have a look at the edge computing architecture. [1]

In edge computing there are mainly three layers:-

2.1. IoT Layer

This layer consists of devices such as sensors, controllers, smart cameras, self-driving automobiles, and cell phones which are known as edge devices. These edge devices are real-time devices which use real-world data.

2.2. Edge Layer

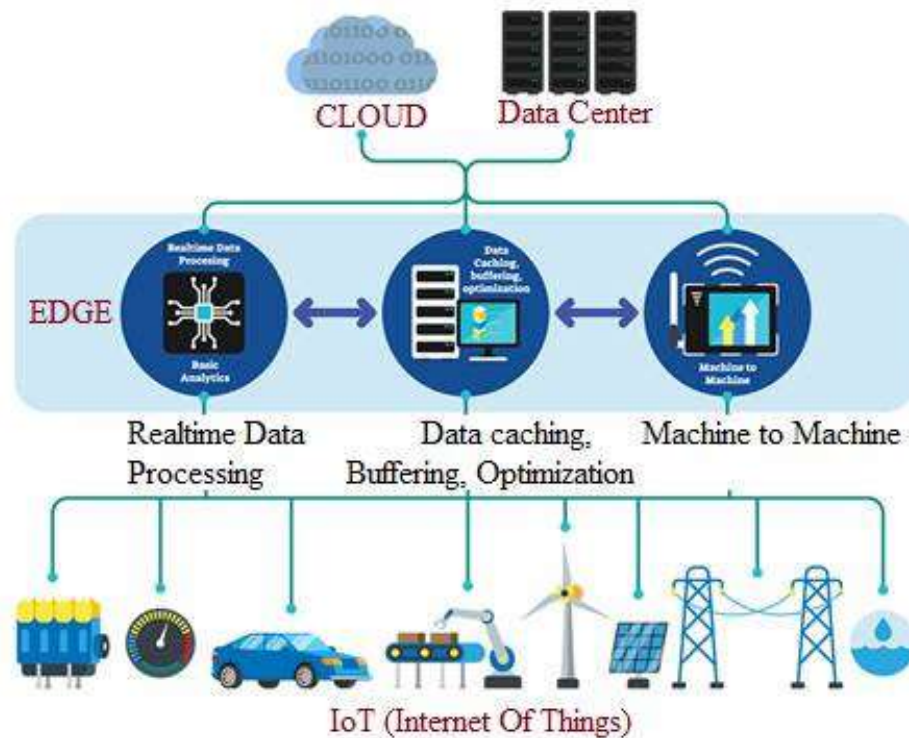


Fig 1.0: Edge Computing Architecture (Image Courtesy: Researchgate.net)

Edge layer is a middle layer which consists of edge nodes. These edge nodes are the computing nodes where data processing takes place. These edge nodes act as a cloudlet which delivers cloud computing services in a rapid manner within a nearby geographical proximity. The edge layer does all the necessary preprocessing required. It results in reduced bandwidth usage, lower latency and better privacy policy enforcement. It is also responsible for basic data caching, buffering, filtering and optimization operations.

2.3. Cloud Layer

Cloud layer is responsible for advanced data processing and warehousing.

3. Edge Computing Implementation

The following decisive steps are required for implementing edge computing.

- Determine the amount of intelligence needed to embed into IOT devices according to the requirement.
- Cluster the various IOT devices in the network in terms of proximity to cut down latency and enhance response time.
- Establish an excellent rules engine for edge server's in tune with the requirements. Also, update the rules engine in accordance with the demands and priorities.
- Keep distinct edge servers in the network along with a central server with optimum intelligence. The data should be dealt with in a hierarchical way in which the data is primarily sent to the least intelligent edge server and subsequently moved to the edge servers with more intelligence up to the central server according to the complexity and priority of data. That means the central

server only processes the critical data.

- Use advanced data analytics in edge servers for precise data analysis, and report generation.[2]

All edge devices are routed through edge nodes where the local processing like caching, basic processing, offloading the data etc... is done. The edge nodes are further connected to a hyper-scale cloud depending on client requirements. [3]

- Behind malls, hospitals, entertainment parks, hotels etc.
- Near telecom towers
- Along railway tracks, ports, airports etc.
- Integrated with community/neighborhood areas
- In factories and industrial parks to support industrial IOT

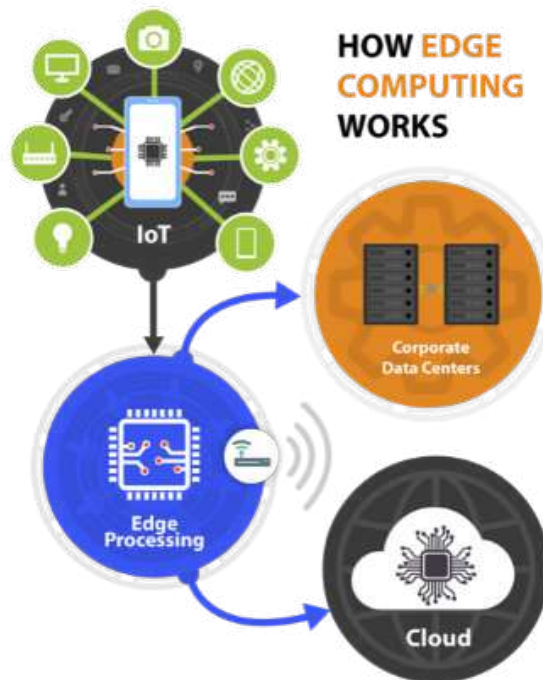


Fig 2.0: Edge Computing Working (Image Courtesy: Spec_india.com)

4. Edge Data Centers

Edge data centers administer end-user computing needs in close proximity to data sources. It may be located primarily at following sites:-

- Buildings/Smart Buildings/Smart Cities
- Co-location within central offices, tech and business parks

- Critical application centre's like disaster recovery, emerging response centers, oil and gas installations, refineries etc.

5. Significance of Edge Computing

5.1. Delivers low latency

Edge computing delivers responses with minimum delay. So it's very helpful in IOT devices where rapid response is required.

5.2. Use high bandwidth

Edge computing rationalizes data downloading and allows each data point to process its own information. It also reduces the storage costs and burden on individual systems. It also boosts browsing.

5.3. Foster real-time automation

Edge computing eliminates the delay in processing time and communication. It improves real-time automation and augments productivity.

5.4. Provides reliability in data storage

Edge computing provides data encryption mechanisms for security which prevents cyber-attacks. The data risk is also reduced as the computing is done in close proximity. It also provides constant data storage. Edge computing sustains the network architecture of the entire distributed system so that fault detection and recovery are easy.

5.5. Optimize equipment up time

Edge computing maximizes the speed and volume of data transfer on a node and provides real-time information on its status and productivity. Also, failures can be easily detected and preventive measures can be taken.

5.6. Provides Data Caching

Here frequently utilized data is cached locally on edge servers which deliver rapid access to data and minimize data traffic between mobile users and the centralized cloud.

5.7. Provides Data Buffering

Here edge devices will be able to buffer the data locally. So that data is not lost even if a network outage occurs. So operations can be continued conscientiously.

5.8. Provides Data Optimization

Edge computing is capable of productively fetching, organizing and analyzing data. In edge computing communication takes place between IOT devices and central Information Technology (IT) networks. [4]

6. Favors of Edge Computing

- The edge can provide latency in milliseconds or even less as data is processed at the data source itself.
- Enhanced throughput.
- Running analytics at the edge helps to reduce the data that has to be sent upstream. It lowers the data costs and extends bandwidth.
- Results in a reduction in data volume to the cloud.
- It scales down data travelling time and transmission cost.
- Cut back connectivity constraints in difficult environments.
- Minimized risk of failure /cyber-attacks/ other vulnerability via the cloud as data is stored, accessed and protected locally.
- Enhanced capacity and device innovation through software updates.

7. Edge Computing Challenges

- Need to preserve reliability and security.
- Lack of visibility into IT & information assets.
- Narrow in-house technical expertise.
- Need to manage massive amounts of disparate devices in the field.
- Due to multiple edge receivers situated at certain distances from the data center, troubleshooting and repair of any issue that arises in the framework need a lot of

logistics as well as manual input which increases the cost of maintenance.

- Drawing huge chunks of data at either edge server or data center the edge nodes may face many security and accessibility issues.
- Various logistic complexities.
- Managing divergent edge network and storage systems are complex and require experienced IT staff available at multiple geographical locations at the same time.
- It bears more financial load on firms.
- As the edge network is small it leads to a lack of power and data computation problems for high data workloads. So various scaling complexities and limited capability issues may occur.[5]

8. Cloud Computing Vs. Edge Computing

See Table 1.0: Cloud Computing vs. Edge Computing [6]

9. Use Cases

9.1. Smart City

- For superlative functioning of smart security devices like cameras, video doorbells etc... for authorization and surveillance. Also IOT sensors, drones, robots etc. With edge computing, these devices can operate with greater security and improved connectivity.
- Automating power grid.
- Assist in the safe operation of autonomous vehicles with various services like real-time communication, image processing, real-time traffic monitoring and mapping, and real-time alerts.
- Assist in the functioning of Smart

Thermostats for automated control of heating, ventilation, and air conditioning remotely.

- Zone priority based on schedule/ event/risk monitoring.
- Can be able to manage an enormous number of active devices or connections at a site.
- Increased technology growth and phased planning.[7]

9.2. Medical

- Edge connected ambulance for safe transportation of patients with live streaming of patients.
- Real-time analysis of medical data like heart rate from sensors, and BP haptic-enabled diagnostic tools for diagnosing ultrasound scanning remotely.
- Efficiently operate various e-health devices and patient wearables like glucose monitors, heart rate, sleep monitoring, BP, body temperature, oxygen saturation, and fall detector.
- Robot-Assisted Surgery.
- Tracking the location of gurney.[8]

9.3. Smart Home

- Assist in automating and intelligent controlling of various heating, lighting, gas, and energy equipment efficiently
- Assist in intelligent operation of Smart Microwave Oven, Smart Kettle with temperature control and automatic scheduling with Google Assistant, Alexa etc...Smart Fan, Smart Refrigerator, Smart Washing machine

9.4. Defence

- Command and control of different combat vehicles like tankers, aircraft, fighter jets etc.

Sl. No.	Factor	Cloud Computing	Edge Computing
1.	Architecture	Centralized. Data centers are remote from client devices.	Distributed and decentralized. Nodes are located in close proximity to client devices.
2.	Communication	Cloud access is more time-consuming.	Less time to access as nodes are close.
3.	Data Processing	Takes place in remote data centers.	Done at the edge of the network.
4.	Computing Capabilities	More powerful with advanced capabilities and storage facilities.	Less powerful.
5.	Latency	High latency.	Low latency.
6.	Security	Less secure.	More secure.
7.	Risk of Failure	Risk of failure is more.	Risk of failure is less.
8.	Response Time	High.	Low.
9.	Bandwidth	Requires huge bandwidth.	Requires low bandwidth.
10.	Technology	Based on the internet-driven global network on robust TCP/IP protocol.	Pushes the intelligence, processing power and communication of an edge gateway or appliance directly into devices.
11.	Use cases and Benefits	Inventory, BI and big data storage, deep analysis, rich data visualization, dashboard reports, back end access	IOT devices, traffic lights, autonomous vehicles

Table 1.0: Cloud Computing vs. Edge Computing [6]

9.5. Finance

- Enhance ATM security by verifying video by integrating image recognition on ATMs. If any security breach occurs it alerts the bank.

9.6. Streaming Services and Content Delivery

- Media companies can use edge computing for uninterrupted video streaming and sharing content.

10. 5G and Edge Computing

5G is the fifth-generation cellular technology with the objectives to enhance speed; curtail latency and boost flexibility. Edge computing is essential for 5G in order to fit the latency targets. It also curtails the bottlenecks in networks and elevates communication. Edge computing also cut down the data processing and transportation costs.

11. Trends in Edge Computing

11.1. Multi Access Edge Computing

Multi Access Edge Computing or Mobile Edge Computing is a network architecture which provides mobile workloads in close proximity to the user using a Radio Access Network (RAN).

It offers low latency, high bandwidth, efficient network operation and service delivery and also enhanced customer experience. Multi Access Edge Computing is mostly used in IOT, location navigation services for mobile devices, data and video analytics etc.

11.2. Fog Computing

In fog computing data is intensely analyzed, and filtered and only important information is stored in the cloud. So that

we can save a lot of space in the cloud and can send data rapidly to the cloud.

12. Conclusion

This paper tried to highlight edge computing architecture and working. Also outlined the importance, and various benefits of edge computing and a brief comparison of edge computing with cloud computing. It also reviewed the various challenges faced by edge computing and some prominent use cases of edge computing. The paper also encompasses a short narration about the role of edge computing in 5G and trends in edge computing. Edge computing is a distributed computing model which accelerates connectivity and brings the data processing in close proximity to the end user.

13. References

<https://www.researchgate.net/figure/Hierarchy-of-Edge-Fog-and-cloud-computing-fig3-342875939/>

https://www.spec_india.com/blog/what-is-edge-computing-the-quick-overview/

<https://www.techtarget.com/iotagenda/tip/How-to-implement-edge-computing-in-5-steps>

[https://www.scitechsociety.com/importtance-of-edge-](https://www.scitechsociety.com/importtance-of-edge-computing-in-industry-4/)

[computing-in-industry-4/](https://www.greyb.com/edge-computing/)

<https://www.greyb.com/edge-computing/>

<https://www.samsolutions.com/blog/fog-computing-vs-cloud-computing-for-iot-projects/>

<https://www.smartcitiesworld.net/opinions/opinions/reviving-smart-cities-with-edge-computing-and-5g/>

<https://www.stlpartners.com/articles/edge-computing/digital-health-at-the-edge/>

Online English Teaching Using Artificial Intelligence Under Big Data Environment

Aparna Suresh, Dr. S.V. Annlin Jeba

*Student, Dept of CSE Sree Buddha College of Engineering, Pattoor, India |
aparnasureshsm@gmail.com*

*Head of the Department, Dept. of CSE Sree Buddha College of Engineering, Pattoor, India |
sureshannlin@gmail.com*

Abstract

Application of big data and artificial intelligence influenced online learning platforms. Here, artificial intelligence and big data are introduced into the English teaching framework to formulate a new and effective teaching Eco-environment for learning. The proposed system meets the needs of social development and international communication in English. In the proposed method, the characteristics of English teaching in a big data environment are analyzed in detail. Then big data technology is used to construct a new Eco-environment of English teaching to improve the teaching and learning quality. The data mining method is used to analyze the relationship of interdependence and mutual restriction among various factors in English teaching in order to build and implement a new Eco-environment with information sharing, quality teaching and personalized learning of English. Finally, the student's adaptability level can be checked by XGB Classifier. Therefore, the constructed Eco-environment provides a new idea and direction for English teaching reform by application of big data and artificial intelligence.

Keywords: *Eco-environment, English teaching, big data, data mining, influence factor, comprehensive ability*

1. Introduction

The rapid development of big data, artificial intelligence, mobile Internet and other modern information technologies, generated a revolutionary reform in the teaching field. The popularity of mobile Internet has created an English learning environment for students, big data technology has opened personalized intelligent teaching, and artificial intelligence technology has promoted the

innovation of English teaching concepts and learning methods. Under the influence of flipped classrooms, massive open online courses, micro classes and other modern teaching technologies, English teaching needs to improve the traditional teaching methods, reverse the relationship between teachers and students, and make teachers change from leaders to guides, while students become the main factor to affect the effectiveness of English teaching. Under the influence of multimedia classrooms and the Internet, the teaching

environment in colleges and universities has changed from a closed environment to an open environment with intelligence, network and digitalization. Therefore, English teachers must consider the reform and promotion of traditional English teaching by using modern information technologies in order to deeply study a new perspective reform on the direction of English teaching. The teaching ecosystem is a relationship between teaching resources and elements. The teaching ecosystem of English regards many factors as interrelated variables, such as students, teachers, textbooks, multimedia, resources, classrooms, schools and society and so on. The traditional English teaching methods cannot meet the ever-changing technology and culture. The new teaching ecology of English must be reconstructed.

2. Related Works

Modern technologies such as big data and artificial intelligence can be used for effective and adaptive learning in the field of online educational platforms. Several works of literature are proposed in this context. This section deals with a detailed discussion of related literature. The paper "A flipped classroom-based teaching system for English teaching," [1] briefs that The flipped classroom (FC) is a novel teaching mode combining traditional classroom and computer network technology. This paper attempts to design an FC-based education system for college English teaching, which is critical to the quality of English teaching in colleges. In our research, a college English education system is established based on SQL server and FC. The architecture and functional modules of the system were discussed in detail. Specifically, the system architecture was created based on the browser/server (B/S) structure, the data management was realized by SQL server 2008, and the system functions were designed and achieved by JavaServer Pages (JSP). The proposed system was tested in an actual application. The results confirm that the system can meet the teaching demand, arouse the student's interest in learning and improve the effect of English teaching. "Timely daily activity recognition from

headmost sensor events,"[2] briefs that Human activity recognition (HAR) has been increasingly used in medical care, behavior analysis, and the entertainment industry to improve the experience of users. Most of the existing works use fixed models to identify various activities. However, they do not adapt well to the dynamic nature of human activities. We investigated activity recognition with postural transition awareness. The inertial sensor data was processed by filters and we used both the time domain and frequency domain of the signals to extract the feature set. For the corresponding posture classification, three feature selection algorithms were considered to select 585 features to obtain the optimal feature subset for the posture classification. And We adopted three classifiers (support vector machine, decision tree, and random forest) for comparative analysis. After experiments, the support vector machine gave better classification results than the other two methods. By using the support vector machine, we could achieve up to 98% accuracy in the Multi-class classification. Finally, the results were verified by probability estimation. The "Teachers' experiences of English-language-taught degree programs within health care sector of Finnish polytechnics,"[3], The purpose of this study was to research teachers' experiences of the English-Language-Taught Degree Programs in the health care sector of Finnish polytechnics. More specifically, the focus was on teachers' experiences with teaching methods and clinical practice. The data were collected from eighteen teachers in six polytechnics through focus group interviews. Content analysis was used to analyze the data. The results suggested that despite the positive interaction between students and teachers, choosing appropriate teaching methods provided a challenge for teachers, due to the cultural diversity of students as well as the use of a foreign language in tuition. Due to students' language-related difficulties, clinical practice was found to be the biggest challenge in the educational process. Staff attitudes were perceived to be significant for students' clinical experience. Further research using stronger designs is needed.

The "Practices that promote English reading for English learners (ELs)," [4], Schools are becoming increasingly diversified; however, training and professional development related to working with English language learners (ELs), especially in the area of English reading, is limited. In this article, we identify three "Big Ideas" of effective and collaborative practices that promote English reading achievement for EL students: (a) foster academic English at all stages of second language acquisition by explicitly teaching vocabulary, emphasizing cross-linguistic transfer strategies, and supporting ongoing oral language development; (b) adopt a schoolwide collaborative approach to conduct frequent formative reading assessments and use the data to drive instruction by providing accommodations that promote English reading; and (c) implement a variety of grouping strategies to deliver reading instruction within a welcoming and sensitive learning climate. In addition, we discuss how school professionals may proactively instruct ELs and collaborate within a multidisciplinary framework to improve the English reading ability of students who are simultaneously learning the English language. "Learning analytics for English language teaching," [5], In recent times, online learning platforms get more and more attention and the number of collected data is growing. Learning analytics is a valuable opportunity to gain specific information for a better understanding of students' learning behavior and to improve their learning success. In this work, the collected data from the online learning platform www.more-online.at is analyzed and the first research results are presented. The main objective is to analyze the usage behavior over a school year and to show the diffusion of the online platform among provinces in Austria, different school types and other characteristics. Furthermore, the content of the online platform is put under closer examination to enable decisions about the efficiency and effectiveness of different types of exercises. "Effect of English corpus on reform of English teaching and the improvement of students' vocabulary

competence," [6], English vocabulary is the basis of language learning. Traditional vocabulary teaching is stereotypical and boring, without sufficient language skills and language knowledge. In this context, the development of corpus provides a platform for English teaching reform, weakening the defects of traditional vocabulary teaching. The authenticity and practicality of corpus can greatly enhance students' interest in learning. For this, this paper studies the effect of the English corpus on college English teaching reform and the improvement of students' vocabulary competence. The research results show that the implementation of English vocabulary teaching reform based on a corpus-assisted platform enhances students' self-learning enthusiasm and increases the internal driving force of students' English learning; through the comparison of the results before and after the corpus teaching experiment, it's found that the in the corpus-based English vocabulary teaching mode, the students' English performance is significantly improved. This study shall provide the theoretical basis for the popularization of corpus in English teaching. "Construction of an ecological teaching model for English courses under the background of Internet plus," [7], College English course is one of the highly practical and compulsory courses. However, the traditional teaching model of current college English is seriously lagging behind modern teaching technology, and the teaching effect is not good. In response to this problem, this paper proposes the use of advanced multimedia technology and network technology in the context of "Internet Plus" to construct an ecological teaching model for college English courses. This teaching model, by making full use of intelligent terminals and wireless campus networks, establishes Internet classrooms, network self-learning centers and language labs. The research results show that the ecological teaching mode realizes the transformation of the traditional teaching classroom to the multimedia and network teaching platform, enhances the modern English teaching system, and thus improves the English teaching effect. "The relationship between English language

learner characteristics and online self-regulation: A structural equation modeling approach," [8] Learner beliefs, anxiety, and motivation are three common learner characteristics. They have consistently been found to account for language learning performance. Meanwhile, self-regulation is critical in sustaining online learners' continuous efforts and predicting their learning outcomes. Despite the massive and rapidly increasing number of online English learners, few studies have clarified the assumed relationships between learner characteristics (learner beliefs, anxiety, motivation) and self-regulation in the online English learning context. This study aims to fill the gap by conducting structural equation modeling analysis to examine their relations. To fulfill the research purpose, we adopted the previous questionnaires with sufficient reliability as instruments to evaluate students' online English learner beliefs, learning anxiety, learning motivation and online self-regulated English learning. The valid responses collected from 425 Chinese undergraduate university students enrolled in an online academic English writing course provided the data source. The results indicated that learner beliefs positively predicted while learning anxiety negatively predicted, online self-regulated English learning. Online English learning motivation was a mediator in these associations. The findings suggested that stronger learner beliefs of self-efficacy and perceived value of English learning promoted learning motivation and self-regulation. In contrast, higher learning anxiety, such as test anxiety and fear of negative evaluation, harmed learners' motivation and their online self-regulated English learning. "Development of a Kinect-based English learning system based on integrating the ARCS model with situated learning," [9] The main design concept was to integrate Kinect as an interaction technique with theories of situated learning and the attention, relevance, confidence, and satisfaction (ARCS) model, to design-relevant learning activities and materials, thereby enhancing students' learning outcomes. The proposed system allows for planning and designing learning activities and content according to situated

learning components and the ARCS model. The somatosensory interaction system Kinect was used to provide users with a virtual learning environment to achieve actual spatial and physical experiences, assisting learners' engagement in stories and scenarios as well as enhancing their motivation to learn. English vocabulary related to supermarkets was set as the learning objective and 70 students ranging from third to sixth grade at a learning center in Tainan, Taiwan were selected as participants. During the experiment, participants were divided into two groups: the experimental group, which employed the proposed learning system, and the control group, in which students learned using printed materials coupled with mobile devices. Pre- and post-test scores of the two groups were used to assess learning outcomes and analyze the ARCS model-based questionnaire. The results revealed that the proposed system effectively improved learners' motivation to learn and learning outcomes. "An enhanced MSIQDE algorithm with novel multiple strategies for global optimization problems," [10] Quantum-inspired differential evolution (QDE) is an evolutionary algorithm, which can effectively solve complex optimization problems. However, sometimes, it easily leads to premature convergence and low search ability and falls to local optima. To overcome these problems, based on the MSIQDE (improved QDE with multi strategies) algorithm, an enhanced MSIQDE algorithm based on mixing multiple strategies, namely, EMMSIQDE is proposed in this article. In the EMMSIQDE, a new differential mutation strategy of a difference vector is proposed to enhance the searchability and descent ability. Then, a new multi-population mutation evolution mechanism is designed to ensure the relative independence of each subpopulation and the population diversity. The feasible solution space transformation strategy is used to achieve the optimal solution by mapping the quantum chromosome from a unit space to a solution space. Finally, some multidimensional unimodal and multimodal functions are selected to demonstrate the optimization

performance of EMMSIQDE. The results demonstrate that the EMMSIQDE is significantly better than the DE, QDE, QGA, and MSIQDE, and has better optimization ability, scalability, efficiency, and stability.

3. Proposed System

The proposed eco-environment applies big data mining for learning requirements, habits, hobbies and so on as the breakthrough point, and combines the textbook writers, teachers, information support institutions, students and other members of the English teaching ecosystem to implement English teaching and learning, model design, resource development, evaluation mechanism and management mechanism. The distributed

file system HDFS and MapReduce model of the Hadoop framework is used to construct the Eco-environment of English teaching with compatible, harmonious coexistence and dynamic development of all sub-systems. The HDFS is mainly responsible for the distributed storage and management of big data, and the MapReduce model is mainly responsible for the calculation and processing of large-scale data. Hadoop uses HDFS to achieve its storage capacity, and MapReduce to achieve its computing capacity. The proposed system consists of four subsystems for Big data analytics. The working of each subsystem in the proposed Eco Environment of English teaching is explained in the following sections:

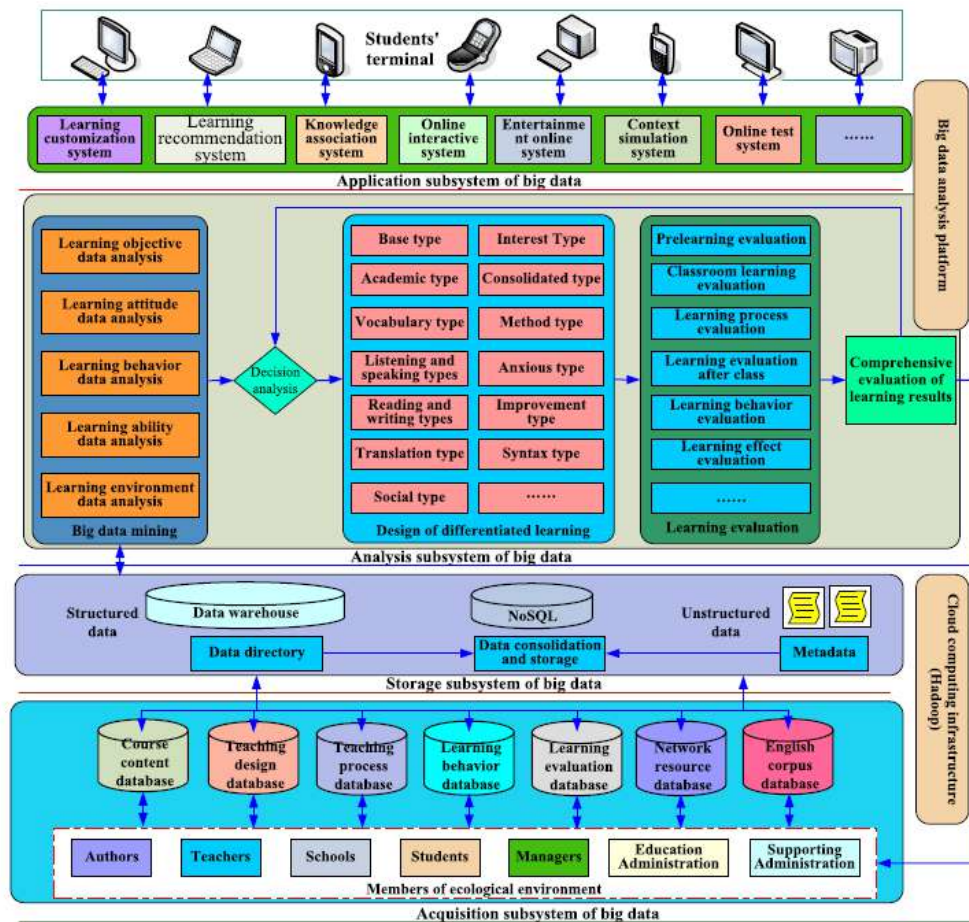


Fig 1: The overall framework of online English learning

Step 1 - The Big Data Acquisition: The working of this step is to collect data from course content, teaching design, teaching process, learning behavior, learning evaluation, network resource and English corpus and so on in order to construct a big data of English teaching for providing a database to analyze.

Step 2 - The Big Data Storage: It preprocesses the acquired big data by cleaning data, archiving and compression, so as to realize the data storage. At the same time, the data of various kinds of English teaching and learning data are processed, summarized and sorted by using data warehouse technology and HDFS distributed file system.

Step 3 - The Big Data Analysis: It achieves the data mining goals by using statistical analysis and artificial intelligence to obtain a dynamic comprehensive evaluation for English teaching, and provide an online learning program.

Step 4 - The Big Data Application: It is to realize these systems of learning customization and recommendation, knowledge association, online interactive, entertainment online and so on by using classification, valuation, forecast, association rule and clustering, which can improve the effect of English teaching and learning.

3.1. Realization of Data Acquisition Subsystem

Data acquisition system is to collect data from data sources to the Eco-environment that can support big data architecture, so as to realize data acquisition and establish a data warehouse in the later stage. Big data acquisition methods mainly include the off-line acquisition method, real-time acquisition method, Internet acquisition method and other data acquisition methods. The acquisition data of English teaching mainly includes textbooks, PPT or multimedia courseware, case teaching, exercise database, test database, real-question database and other content teaching materials. The acquisition data includes the compilation of the teaching syllabus, teaching purpose, learning

content and so on in the teaching design, the emotional attitude, context construction, and so on in the teaching process, the self-test, learning time, approach, motivation and attitude and so on in the learning behavior, the evaluation data of learning strategy, style, ability, process, result and so on in the learning evaluation, the English teaching and learning materials in the network resources, the English textbooks, English auxiliary materials, English and American novels, essays, scripts, press releases and other corpus in the English other corpora Eco-environment of English teaching, the acquisition data mainly includes structured data and unstructured data, such as words, numbers, graphics, images, videos, animations, audio, and so on, which are used to realize data acquisition and sharing. The big data acquisition subsystem integrates various structured data and unstructured data to provide a database for the efficient use of teaching data resources. Therefore, in order to realize the big data acquisition subsystem, we developed an adaptive interface. The corresponding interface module is developed to interface with various information systems for the existing information system. The realized subsystem can realize the data-sharing interface. The SQL server 2018 is selected to realize the unified storage and management of data. We developed the relevant interface to obtain the relevant data information according to the data situation in order to complete the data acquisition and extraction.

3.2. Realization of Data Storage Subsystem

In the Eco-environment of English teaching, the big data storage subsystem mainly preprocesses the acquired massive structured data and unstructured data, including data cleaning, archiving and compression, so as to realize the integration of data storage. The deep integration of English teaching data can solve the problem of information islands caused by different countries, regions and schools, as well as different data storage structures and operating systems. The data warehouse is used to support the

management decision-making process. It is a general term for integrating operational data into a unified environment in order to provide decision-making data access.

The data warehouse of English teaching is organized according to a certain theme. It is to process, summarize and sort the data of various kinds of original and scattered data of English ecological teaching and learning, in order to eliminate the inconsistency between various kinds of English ecological teaching and learning data. Therefore, the data warehouse of English teaching is the consistent data of teachers, students, institutions and other overall situations. These data are stored in the cloud database through a cloud platform, which can provide data support for the analysis and application of big data. HDFS is an effective tool to realize distributed storage and management of large-scale data. It adopts a typical master/slave structure, which greatly simplifies the system architecture, and makes the system to be more concise and convenient for system management. The data node in the file system mainly stores the actual data and is mainly responsible for the storage management on the physical node. From the internal structure of the distributed file system, we divided data files into multiple data blocks, which are stored in each data node. Each data node stores data blocks from multiple files. Therefore, the HDFS of Hadoop is used to store all kinds of collected data in a unified method to improve the scalability and fault tolerance of data storage. According to the corresponding rules, the acquired data is stored as a set of complete data file set and forms a data warehouse.

3.3. Realization of Data Mining Subsystem

In the Eco-environment of English teaching, data mining is the process of searching and hiding information from a large amount of data by algorithms. We use statistics, information retrieval, machine learning, expert systems and pattern recognition to achieve the goals of data mining. In the Eco environment of English teaching, the regression analysis is

used to analyze the individual English learning behaviour factors according to the English learning objectives. The learning motivation, attitude, interest, perseverance, ability, method, habit and other factors of students are analyzed in detail by using a statistical analysis method. The mined results are classified, valued, predicted and clustered by using clustering algorithm, association rules, neural network algorithm and so on. According to the learning behaviours and attitudes of students, these students will be divided in order to determine the similar learning requirements of different students, as well as the differences in learning requirements of students in different learning stages. At the same time, the big data mining subsystem will make a dynamic comprehensive evaluation of the teaching process, learning process and learning results, and determine the advantages and disadvantages of the teaching program according to the learning evaluation results. Here we are using XGBClassifier to check the student's ability. XGBoost classifier is a machine learning algorithm that is applied for structured and tabular data. XGBoost is an implementation of gradient-boosted decision trees designed for speed and performance. Finally, we adjusted the online and off-line learning programs.

3.4. Realization of Data Application Subsystem

The big data application subsystem mainly includes the learning customization system, learning recommendation system, knowledge association system, online interactive system, entertainment online system, context simulation system, online test system, and so on. By connecting the learning terminal to the relevant application system, these students can learn English. The learning customization system meets the requirements of these students to customize the learning information according to their own learning demands. The learning recommendation system automatically recommends learning information to the students. The knowledge association system automatically associates and recommends the corresponding extended

knowledge. The online interactive system communicates with teachers and solves these problems in the process of English learning. At the same time, it can communicate with other students and share learning experiences. The context simulation system realizes the online simulation of English learning circumstances and lets students practice English communication in the English language circumstances closer to the real environment.

Algorithm: English Learning Eco-environment

Input: Course and interactive learning parameters

Output: An effective English Learning Environment

START

Step 1: Data acquisition

1.1. Data is collected from members of the ecological environment and stored in different databases

Step 2: Data Storage

2.1 Storage of structured and unstructured data into data consolidation and storage space

Step 3: Big Data Analysis

3.1 Big data mining of learning attributes

3.2 Design of Differentiated Learning

3.3 Evaluate learning effectiveness in the new environment

Step 4: Application of mined big data results using Artificial Intelligence

13. References

H. Y. Zhang, "A flipped classroom-based teaching system for English teaching," *Int. J. Emerg. Technol. Learn.*, vol. 14, no. 16, pp. 120–129, 2019.

Y. Liu, X. Wang, Z. Zhai, R. Chen, B. Zhang, and Y. Jiang,

4.1 Develop a learning customization system using XGB Classifier

STOP

XGBoost is an implementation of Gradient Boosted decision trees. In this algorithm, decision trees are created in sequential form. Weights play an important role in XGBoost. Weights are assigned to all the independent variables which are then fed into the decision tree which predicts results. The weight of variables predicted wrong by the tree is increased and these variables are then fed to the second decision tree. These individual classifiers/predictors then ensemble to give a strong and more precise model. It can work on regression, classification, ranking, and user-defined prediction problems. Thus, it will provide an accurate prediction of a student's ability.

4. Conclusion

The application and development of big data technology and artificial intelligence methods have changed the balance among the factors in the traditional English teaching ecosystem. At the same time, the existence of big data technology and artificial intelligence methods promotes the construction of a new Eco-environment for English teaching. An Eco environment of English teaching is constructed with interdependence, mutual restriction and rebalancing of all factors by using the distributed file system HDFS and MapReduce model of the Hadoop framework. A new teaching mode and method of English based on combining big data and artificial intelligence is obtained. It has an important guiding role in improving the teaching quality of English teaching. Student's adaptability is checked by the XGB Classifier.

"Timely daily activity recognition from headmost sensor events," ISA Trans., vol. 94, pp. 379–390, Nov. 2019.

M. Pitkajarvi, E. Eriksson, and P. Kekki, "Teachers' experiences of English language taught degree programs within health care sector of finnish

polytechnics," Nurse Edu. Today, vol. 31, no. 6, pp. 553–557, Aug. 2011

R. S. Martínez, B. Harris, and M. B. McClain, "Practices that promote English reading for English learners (ELs)," J. Educ. Psychol. Consultation, vol. 24, no. 2, pp. 128–148, Apr. 2014.

H. Volk, K. Kellner, and D. Wohlhart, "Learning analytics for English language teaching," J. Universal Comput. Sci., vol. 21, no. 1, pp. 156–174, 2015.

S. Wang and X. Zeng, "Effect of English corpus on reform of English teaching and the improvement of students' vocabulary competence," Teach. Sci.-Theory Pract., vol. 18, no. 6, pp. 3493–3499, 2018.

B. Wu, "Construction of ecological teaching model for English course under the background of Internet plus,"

Teach. Sci.-Theory Pract., vol. 18, no. 6, pp. 3515–3521, 2018.

W. Wang and J. Zhan, "The relationship between english language learner characteristics and online self-regulation: A structural equation modeling approach," Sustainability, vol. 12, no. 7, p. 3009, Apr. 2020.

Y.H. Chang, P.-R. Lin, and Y.-T. Lu, "Development of a Kinect based english learning system based on integrating the ARCS model with situated learning," Sustainability, vol. 12, no. 5, p. 2037, Mar. 2020.

W. Deng H. Liu, J. Xu, H. Zhao, Y. Song, "An enhanced MSIQDE algorithm with novel multiple strategies for global optimization problems," IEEE Trans. Syst., Man, Cybern. Syst., to be published, doi: 10.1109/TSMC.2020.3030792

Employee Attrition Prediction Using Various Machine Learning Algorithms

Anjaly Denny K., Resija P.R.

*Second year MSc Computer Science, Vimala College (Autonomous), Thrissur |
anjalydennyk@gmail.com*

*Assistant Professor, Department of Computer Science, Vimala College (Autonomous),
Thrissur | resijapr1995@gmail.com*

Abstract

We discuss the tremendous increase in the employee attenuation rate in many organizations. It has been observed to be more severe in small-scale to large-scale industries. Loss of talented employees from companies is a major problem encountered by business leaders. This paper studies employee attenuation using different machine learning models. The model predicts whether an employee is about to leave or not from the institution in the near future and also predicts the reasons in search for a different work environment. It will help them to recognize the employees who are not interested in the work, monthly income, overtime, environment satisfaction, distance from home etc. For analyzing employee attrition we can implement different classification methods such as Logistic Regression, Decision Tree, Random Forest, SVM, KNN, Naive Bayes, and Ada boost Classifier methods on the available dataset. In this paper, we also find out which classification model is more accurate to find out the employee attenuation rate from various institutions.

Keywords: Attrition, Ada boost Classifier

1. Introduction

A huge increase in employee attrition rate and its impact has been observed in IT companies. In this article, we are discussing employee attrition prediction, that is, predicting that an employee will leave the current company (or will resign from the current company) and we will do this using several machine learning algorithms. There are many possible reasons why an employee decides to leave the company. The employee decides to leave an organization because of one of the following reasons or more.

- More salary or good job opportunities
- Not satisfactory working environment
- Huge workload stress
- For higher studies
- Not feeling appreciated by colleagues
- Seeing well-trained employees leave the company

Here we take the dataset given by IBM. The

dataset consists of 15 features containing categorical and numeric features. From this dataset, we perform data pre-processing and classification techniques on the test dataset to get the necessary results which are necessary to stop employee attrition.

1.1. Solution Approach

Many companies are trying to find out various ways to retain their workforce like generating yearly people survey reports and finding concerns of employees. In this paper, we are applying a Machine Learning algorithm to the data collected from the IBM dataset. This can be used by the HR team to identify the reasons behind attrition which are specific to their firm and will give an early indication and time to retain before it is late. Need of Employee Attrition Prediction.

- If HR of one particular project came to know about the employee who is willing to leave the Managing workforce: If the supervisors or HR came to know about some employees are planning to leave the company then they could get in touch with those employees so that we can help them to stay back or by hiring the new alternative of those employees we can manage the workforce.
- Smooth pipeline: If all the employees in the current project are working continuously then the pipeline of that project will be smooth but the workflow will not be so smooth if one employee suddenly leaves that company.

2. Related Work

Voluntary employee attrition is one of the major concerns for any company due to the severity of its impact. Talented employees are a major factor in business success and replacing such talent can be difficult and time-consuming [1]. Researchers have studied that several factors can strongly contribute to employee attrition. Money is not the only factor, as other combinations of factors, such as workload, performance pay and a

weak career plan, have increased the attrition rate in the retail industry. Several studies have explored the use of machine learning to predict employee behavior. In [2], the authors explain in detail how to address employee attrition problems across IT firms using the machine learning model. This model predicts whether a resource is about to leave or not in the near future and predicts the possible reasons behind leaving. The employee decides to leave an organization because of many reasons such as more salary or a better role offered outside, not a suitable working environment, not feeling aligned with company goals, Work-life balance missing, stress, not feeling appreciated or feeling underutilized, and seeing good employees leave. Machine learning algorithms are applied to data collected from both Alumni and Current employees of a company and create a model. Therefore we can know the reasons behind attrition. This paper explains the basic pipeline of predicting workforce attrition in IT companies. It included basic exploratory data analysis, feature engineering, implementing Logistic Regression, Random Forest, Decision Tree and Gradient Boosting. Features like Travel time, Involvement in the job, Working environment, Job location, Annual rewards, Travel opportunities, Learning opportunities, Years with the manager, and Years since the last promotion were common in top ranks in all models. In [3], the authors explored several machine learning algorithms to predict employee churn (or attrition) such as Logistic regression, Naive Bayes, Decision Tree, Support Vector Machine (SVM), and K-Nearest Neighbour (KNN). For that, they take employee data from IBM which contains 1470 records and 35 fields including categorical and numeric features. Employee attrition identification helps in predicting and resolving the issues of attrition. The results of this learning describe that data extraction algorithms can be utilized to construct reliable and accurate predictive methods for employee attrition. The issue of attrition identification is not just to depict attritioners from non attritioners. By using tentative data study and data extraction

methods we can depict the attrition probability for each employee. Authors in [4] categorized employee attrition into two categories voluntary and involuntary attrition. Employee attrition can be defined as the loss of employees due to any of the following reasons such as personal reasons, low job satisfaction, low salary and bad business environment. Involuntary attrition occurs when employees are terminated by their employer for different reasons, such as low employee performance or business requirements. In voluntary attrition, on the other hand, high-performing employees decide to leave the company of their own volition despite the company's attempt to retain them. All trained models are evaluated by measuring their accuracy, precision, recall and f1 score. In this research Random Forest, SVM and KNN classification models are evaluated. It was found that the top three features were overtime, total working years and job level. The main objective of this research was to use machine learning models to predict employee attrition based on their features. As a result, this will help management to act faster to reduce the likelihood of talented employees leaving their company. Previous studies presented different accuracy measures where they used different machine learning models and various data sets. As a result, it is difficult to conclude which model is the best to use.

3. Methodology

3.1. Architecture

Model consists of an employee data set and different machine learning techniques and data pre-processing techniques. In this data set, we have different features such as Location, Emp. Group, Function, Gender, Tenure, Tenure Grp., Experience (YY.MM), Marital Status, Age in YY., Hiring Source, Promoted/Non-Promoted, Job Role Match etc. Based on these features we reach a conclusion about whether the employee will leave the company or not. The above-mentioned features became the predictor variables of the dataset. The target variable is "Stay/Left" with the values

Stay and Left. we have to import various libraries such as Pandas, Numpy, Matplotlib for the prediction. We have to check whether there are null values in the data set. We can see there are null values in the Experience and Job roles so we have to drop them.

3.2. Model Development

For predicting employee attrition we use some methodologies of data classification.

3.3. Logistic Regression

It is used for solving the binary classification problem as well as it is a supervised regression method. We can compute the probability of an event occurring. Logistic Regression has a statistical approach for assessing a data set in which there are one or more autonomous variables that establish an outcome [5]. It is used for predicting the categorical dependent variable using a given set of independent variables.

3.4. Decision Tree

A Decision Tree is a tree-like structure that consists of branches, root nodes and leaf nodes [6]. Each internal node denotes a test on an attribute, each branch shows the result of a test, and every leaf node represents a class label. It is a graphical representation for getting all the possible solutions to a problem or a decision based on the given condition. A decision tree simply asks a question (Yes/No). It further splits the tree into subtrees. Decision tree classification algorithm used for categorical data working of decision tree algorithm for predicting the class of a given data set the algorithm starts from the root node and compares the value of root attribute with the record attribute and based on the comparison follows the branch and jump to the next node process continues and it continues until it reaches the leaf node of a tree.

A. Steps

Step 1: Begin the tree with the root node,

says S which contains a complete data set.

Step 2: Find the best attribute in the data set using ASM (Attribute selection measure)

Step 3: Divide the S into subsets that contain possible values for the best attributes

Step 4: Generate the decision tree node which contains the best attribute

Step 5: Recursively make a new decision tree using the subsets of the data sets created in step 3. Continue this process until a stage is reached where you cannot further classify the nodes that is, leaf node

3.5. Random Forest

It can be used for both classification as well as regression and it is a supervised machine learning algorithm [7]. Decision trees are created on randomly selected data samples. It makes a random selection of features rather than using all features to develop trees. It is a very good indicator of feature importance and it is considered as a highly accurate and robust method because of the number of decision trees involved in the process. Random Forest builds the forest with various decision trees. A collection of numerous random trees called random forests.

A. Working of Random forest algorithm

Step 1: Select random K data points from the training set

Step 2: Build the decision tree associated with the selected data points

Step 3: Choose the number n for the decision trees that we want to build

Step 4: Repeat step 1 and 2

Step 5: For the new data points find the predictions of each decision tree and assign new data points to the category that wins the majority votes.

3.6. KNN Classifier

It is the simplest machine learning algorithm. It is most suitable for classification problems. The KNN algorithm classifies a new data point based on similarity. When a new category appears we can classify them based on their similarity, K-Nearest Neighbour compare each value with the neighbour value [8]. The KNN algorithm assumes the similarity between the new case/data and available cases and puts the new case into the category that is most similar to the available categories.

A. Steps of KNN algorithm

Step 1: Select the number of K neighbors

Step 2: Calculate the Euclidean distance of K number of neighbors

Step 3: Take the K nearest neighbors as per the calculated Euclidean distance

Step 4: Among the K neighbors count the number of data points in each category

Step 5: Assign the new data points to that category for which the number of neighbors is maximum

Step 6: Model is ready.

3.7. Support Vector Machine

Supervised Learning algorithms, which is used for Classification as well as Regression problems. It is used for Classification problems in Machine Learning. We can create the best line or decision boundary that can segregate n-dimensional space into classes. It is the main goal of the SVM algorithm. We can easily classify new data points in the correct category in future. This best decision boundary is called a hyperplane. SVM finds out the vectors that help in creating the hyperplane. Support Vector Machine classification is based on identifying the hyperplane that entirely differentiates the vector into two non-overlapping classes [9]. These vectors

which describe the hyperplane are the support vectors and therefore it is called Support Vector Machine.

3.8. Naive Bayes Classifier

Naive Bayes is based on the Bayes theorem and it is a classification algorithm [10]. It is a family of algorithms, not a single algorithm and they all share a common principle which means that every pair of features being classified is independent of each other. Naive Bayes classifier shows the existence of a specific in a class is unrelated to the existence of any other feature. It is mainly used for test analysis and sentiment analysis.

Naive: It is called Naive because it assumes that the occurrence of a certain feature is independent of the occurrence of the other feature

Bayes: it is called Bayes because it is based on Bayes theorem. Bayes rule is used to determine the probability of a hypothesis with prior knowledge. It depends on conditional probability.

A. Steps of Naive Bayes algorithm

Step 1: Convert the given data set into a frequency table

Step 2: Generate a likelihood table by finding the probabilities of a given feature

Step 3: Use Bayes theorem to calculate the posterior probability

3.9. AdaBoost Classifier

AdaBoost is an ensemble learning method which was initially created to increase the efficiency of binary classifiers. AdaBoost uses an iterative approach to learn from the mistakes of weak classifiers and turn them into strong ones

4. Findings

Exploratory Data Analysis Here we take each feature in the dataset and analyze how they affect the attrition rate in the

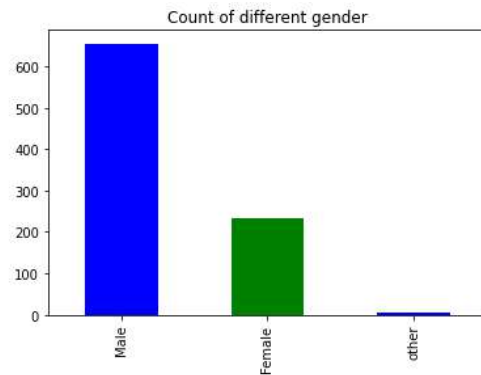


Fig. 1: Count of different gender

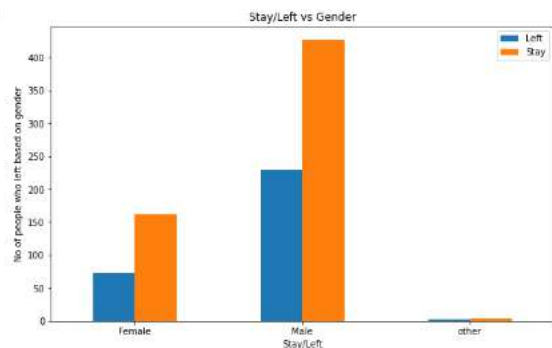


Fig. 2: Gender v/s Stay/Left

company.

From the above graph, we can observe that the count of males is more than other categories of gender. So here male employees have more chances to leave the organization.

Now we can analyze how Gender becomes a feature for the attrition of employees from the company. Here, from the graph we can identify that it heavily depends on males, also we can see that it's either male, female or others but most of them are staying in the company. In this manner, we can also find out how the attributes such as Promoted/Non-Promoted, Function, Job Role Match, Age, and Experience depend on employee attrition from the company.

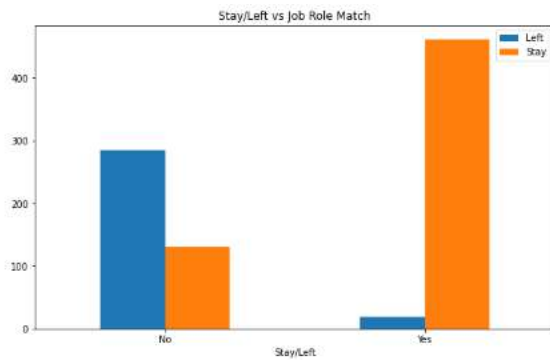


Fig. 3: Stay/Left vs Job Role Match

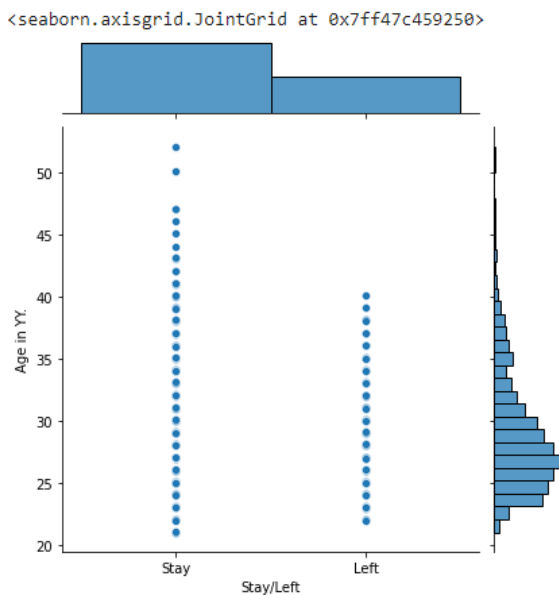


Fig. 4: Stay/Left vs Age

From the graph it is visible that the employees who are not promoted are leaving the company more as compared to the ones who are promoted. We can analyze that the maximum number of employees are in the operation section and a high number of employees in the same section are staying in the company. We can analyze that the number of employees who got the correct job role is staying in the company rather than the ones who don't have the right job role. Employees who are having more age are

staying back in the company rather than the ones who have comparatively less age. Employees who have more experience will stay back in the company rather than the ones who have comparatively less experience. These graphs clearly show the reasons behind the attrition. A person who doesn't have much experience and who doesn't have the right job role is not ready to withstand the company.

4.1. Data cleaning

After values are encoded, build a new location to be used to categorize data columns. Here, we are using integer values instead of the actual region name in location dict new so that our machine learning model could interpret it. We will make a function for the location column to make a new column because our machine learning algorithm will only understand int/float values.

4.2. Correlation matrix

Here, first we are trying to get the correlation between variables where the dataset is not processed that is why we are not able to see the results in the manner we want to, but in later we will see the better correlation plot with the help of processed data

4.3. get dummies ()

For manipulating data, the Pandas get dummies () function is used. Using this function we can convert the categorical values to dummy variables and this has been done with

- Function
- Hiring Source
- New Marital
- New Gender
- Tenure group

values are converted into 1 and 0 respectively for encoding purposes. We are

Machine learning Models	Accuracy
Logistic regression	0.877095
KNN Classifier	0.586592
Support Vector Machine	0.865922
Random Forest Classifier	0.893855
Decision Tree	0.849162
Naïve Bayes Classifier	0.849162
Ada boost Classifier	0.888268

Fig. 5: Accuracy Table

using the Marital () function to convert categorical values Yes and No into integer values that is, 1/0 so that we are assigning the New Marital So now we have to concatenate the columns which are being cleaned, sorted, and manipulated by us as processed data. Then we have to drop the unnecessary columns Then we get the correlation plot of the processed data. Now we have to save the cleaned dataset into another CSV file. Then separate the features and target column again. We have to develop various machine learning models to find out the attrition rate. By comparing each model we find out which model is the more accurate one among them. So here we use seven different machine learning algorithms such as Decision Tree, Logistic Regression, Naive Bayes, Random Forest, K Nearest Neighbour, Support Vector Machine, and Ada boost Classifier. We find a classification report for each model. Here we calculated precision, Recall, f1-score, and support of every model.

5. Simulations and Result

Finally we compare each model with all other models. So that we can find out the most accurate model. Here, we can see that AdaBoost Classifiers have the best accuracy..

6. Conclusion

Employee attrition is a major issue. It is difficult to substitute a well-trained employee when there is attrition and it is cost-effective. We try to find out and analyze the employee information to estimate future attritioners and study the reasons for employee turnover. Machine learning algorithms can be used to construct reliable and accurate predictive models for employee attrition. Here we find out the precision, accuracy, and f1 score of each machine learning model and find out which model is more accurate by comparing other ones. Finally, we reached

```
#Visualize the accuracy of each model
model_compare.T.plot(kind='bar') # (T is here for transpose)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7ff479960a90>
```

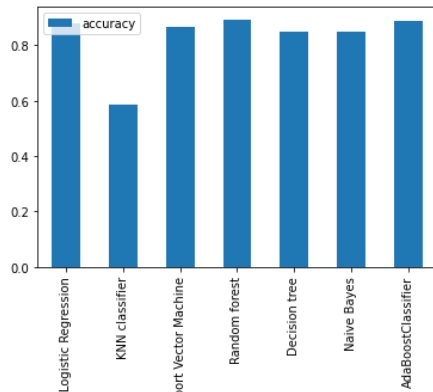


Fig. 6: Visualizing the Accuracy Result

the conclusion that the AdaBoost Classifier algorithm is the most accurate one compared to other machine learning algorithms.

7. Acknowledgement

We would like to express our sincere

gratitude to all who helped to complete the work. We are also thankful to our friends and family members for their valuable support and encouragement. Moreover, we are thankful to God Almighty for showering his choicest of blessings throughout our way.

8. References

S. Kaur and M. R. Vijay, "Job satisfaction-a major factor behind attrition of retention in retail industry," *Imperial Journal of Interdisciplinary Research*, vol. 2, no. 8, pp. 993-996, 2016.

P. Sadana and D. Munnuru, "Machine learning model to predict workforce attrition," in *Proceedings of the 2nd International Conference on Recent Trends in Machine Learning, IoT, Smart Cities and Applications*, pp. 361-376, Springer, 2022.

R. S. Shankar, J. Rajanikanth, V. Sivaramaraju, and K. Murthy, "Prediction of employee attrition using datamining," in *2018 IEEE International Conference on System, Computation, Automation and Networking (ICSCAN)*, pp. 1-8, IEEE, 2018.

S. S. Alduayj and K. Rajpoot, "Predicting employee attrition using machine learning," in *2018 International Conference on Innovations in Information Technology (IIIT)*, pp. 93-98, IEEE, 2018.

X. Zou, Y. Hu, Z. Tian, and K. Shen, "Logistic regression model optimization and case analysis," in *2019 IEEE 7th International Conference on Computer Science and*

Network Technology (ICCSNT), pp. 135-139, 2019.

A. Navada, A. N. Ansari, S. Patil, and B. A. Sonkamble, "Overview of use of decision tree algorithms in machine learning," in *2011 IEEE Control and System Graduate Research Colloquium*, pp. 37-42, 2011.

I. Arel, D. C. Rose, and T. P. Karnowski, "Deep machine learning-a new frontier in artificial intelligence research [research frontier]," *IEEE computational intelligence magazine*, vol. 5, no. 4, pp. 13-18, 2010.

K. Taunk, S. De, S. Verma, and A. Swetapadma, "A brief review of nearest neighbor algorithm for learning and classification," in *2019 International Conference on Intelligent Computing and Control Systems (ICCS)*, pp. 1255-1260, 2019.

M. Hearst, S. Dumais, E. Osuna, J. Platt, and B. Scholkopf, "Support vector machines," *IEEE Intelligent Systems and their Applications*, vol. 13, no. 4, pp. 18-28, 1998.

S. Angra and S. Ahuja, "Machine learning and its applications: A review," in *2017 International Conference on Big Data Analytics and Computational Intelligence (ICBDAC)*, pp. 57-60, IEEE, 2017.



An Authenticatable (2, 3) Secret Sharing Scheme Using Meaningful Share Images Based on Hybrid Fractal Matrix

Veena P., Dr. Anil A.R.

Department of computer science engineering, Sree Buddha College of Engineering Pattoor, Alappuzha

Associate Professor, Department of computer science engineering, Sree Buddha College of Engineering Pattoor, Alappuzha

Abstract

Secret image sharing (SIS) allows a user to create the share images from a secret image in such a way that an individual share does not reveal any information about the secret image, however, when a specified number of shares are brought together, they can be used to reconstruct the secret image. Secret image-sharing mechanisms have been widely applied to the military, e-commerce, and communications fields. This project aims to describe the design and development of a basic Secret Image Sharing System that is capable of efficiently generating shares and reconstructing the secret image from the shares. This system provides a framework which can be used as a powerful tool for Secret Image Sharing research. Efficient calculations of data make the framework an extremely powerful one. It distributes the secret data into multiple shares so that the participants can obtain the embedded secret by sharing their authenticated shares. Through a simple process of these shares, the secret data can be extracted. In this paper, we propose a (2, 3) SIS scheme based on a fractal matrix. Through the guidance of the proposed fractal matrix, the secret data can be distributed into three shares which are indistinguishable from their corresponding cover images. Any two of the three distinct shares can cooperate to extract the exact secret data. Moreover, we devise two authentication mechanisms to prevent tampering. Experimental results show that our proposed scheme can provide efficient payloads with shares of good visual quality. The authentications are also effective and easy to implement.

Keywords: *Secret image sharing, fractal matrix, data hiding, steganography, authentication.*

1. Introduction

Handling secrets has been an issue of prominence from the time human beings started to live together. Important things and messages have always been there to

be preserved and protected from possible misuse or loss. Sometimes a secret is thought to be secure in a single hand and at other times it is thought to be secure when shared in many hands. Some of the formulae of vital combinations of

medicinal plants or roots or leaves, in Ayurveda were known to a single person in a family. When he becomes old enough, he would rather share the secret formula to a chosen person from the family, or from among his disciples. There were times when the person with the secret died before he could share the secret. Probably, similar incidents might have made the genius of those eras think of sharing the secrets with more than one person so that in the event of death of the present custodian, there will be at least one other person who knows the secret. Secret sharing in other forms were prevailing in the past, for other reasons also. Secrets were divided into a number of pieces and given to the same number of people. To ensure unity among the participating people, the head of the family would share the information with respect to wealth among his children and insist that after his death, they all should join together to inherit the wealth.

The rapid development of the internet allows people to continuously enjoy the fruits of the latest technological revolution. More and more information is transmitted through the Internet. However, network viruses, Trojan horses, and illegal organizations that steal information all pose serious threats to information security. Thus, the security of information transmission has become particularly important, especially in the political, military, and commercial fields. Although traditional encryption can ensure information security to a certain extent, this type of technology does not conceal the content of the information. On one hand, the secret keys generated in order to hide information are often unordered symbols, which will arouse the interest of the tracker, thereby increasing the risk of exposure. This has become the biggest weakness of traditional encryption algorithms. On the other hand, with the rising sharp use of digital media, protecting intellectual property rights in digital media, preventing illegal copying and dissemination of intellectual property products, and ensuring information security has also become more and more urgent. In order to solve the above

mentioned series of problems, information hiding [1] technology was born.

Information hiding is mainly divided into two categories: reversible information hiding [2], [3] and irreversible information hiding [4], [5]. In reversible information hiding, the recipient can restore the original carrier losslessly after extracting the secret data embedded by the sender, while after irreversible information hiding, one cannot restore the original carrier. Based on these two categories of information hiding methods, it can be subdivided into the following four types: information hiding based on spatial domain [6]–[12], information hiding based on frequency domain [13]–[15], information hiding based on encryption domain [16]–[19], and information hiding based on compressed domain [20], [21].

Traditional information hiding schemes have an obvious common disadvantage: only a single carrier is used to embed information. Once the carrier is targeted by an attacker, the possibility of the information being attacked increases greatly. To solve this shortcoming, Shamir [22] and Blakley [23] proposed a (k, n) -threshold secret sharing scheme in 1979. Their scheme pioneered the concept of secret sharing. It distributes the secret data into n parts and deals them to n participants. When extracting the secret data

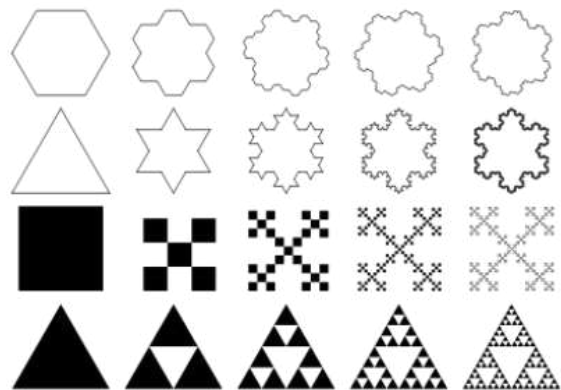


Fig.1: Fractal image

($k < n$) participants working together can obtain the complete secret data. In 1994, Naor and Shamir [24] applied the concept of secret sharing to images and proposed the visual secret sharing (VSS) scheme. But their scheme produces binary shares and these shares look meaningless. A meaningless share is easily noticed by an attacker during transmission. Covering the secret data with a regular image can effectively reduce the risk of being suspected. To solve these disadvantages of the original VSS scheme, many improved VSS methods [25]–[28] have been proposed.

In 2006, Chang et al. [29] proposed a reversible secret image sharing (SIS) scheme using the exploiting modification direction (EMD) reference matrix. After embedding, two steganographic images with high visual quality are generated. The main difference between VSS and SIS is visual perceptibility. Traditional VSS (including visual cryptography [25], extended visual cryptography [26], and color extended visual cryptography [27], [28]) reconstruct images by mixing gray levels or colors through stacking. The SIS scheme proposed in [29] extracts the secret data by computation. In this modern age of consumer electronics, portable devices are popular. Therefore, thin client computing is now more feasible than stacking transparencies. The SIS scheme became a new hot issue.

In 2010, Chang et al. [30] proposed a reversible SIS scheme based on the Sudoku matrix and Lagrange polynomial. In 2018 and 2019, reversible and authenticatable SIS schemes were proposed based on the turtle shell matrix [31], [32]. In 2020, Gao et al. [33] designed a new reference matrix called the stick insect matrix to further improve the SIS schemes. These reversible SIS schemes use only a single cover image to generate multiple shares, where each share has a very slight and invisible difference.

Although the shares generated from a single cover image look the same, they are vulnerable to steganalysis. By analyzing the differences between shares,

camouflage can easily be detected. In 2019 and 2020, SIS schemes using different cover images were proposed [34], [35]. However, these schemes are irreversible: the cover images cannot be reconstructed after the extraction of the secret data.

With irreversibility comes the problem that when one of the shares is tampered with, the secret data cannot be recovered at all, even though the share is authenticatable. In order to reduce the risk of tampering attacks, Gao et al. [36] proposed a (2, 3) reversible SIS scheme based on the fractal matrix. Through the ingenious design of the fractal matrix, any two out of the three shares can cooperate to retrieve the secret data. However, this scheme uses a single cover image to generate two shares and suffers from the steganalysis risk mentioned above.

In this paper, we propose a (2, 3) SIS scheme which is based on the new fractal reference matrix and which uses three distinct share images. In our SIS scheme, three shares that are indistinguishable from their corresponding cover images are generated and any two shares can cooperate to retrieve the secret data. Our scheme also provides two effective authentication mechanisms and a large payload. The experimental results show that the embedding capacity of our scheme is larger than those of related works and that the visual quality of the shares is satisfactory.

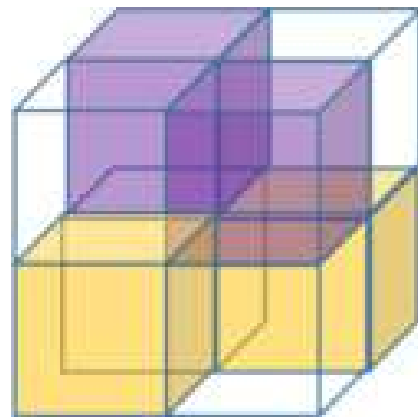


Fig. 2(a): A fractal model

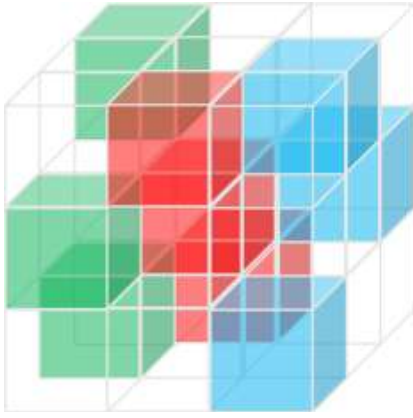


Fig. 2(b): A fractal model

Fig.2: Illustration of the fractal models

The rest of this paper is as follows. Section 2 briefly introduces the new proposed fractal matrix. The share generation, authentication, and data extraction processes are proposed in Section 3. Section 4 presents the experimental results and discussions. Finally, conclusions and prospects are presented in Section 5.

2. The Proposed Fractal Matrix

In this section, we briefly introduce the fractal matrix proposed by Gao et al. in [36]. Then, a modified fractal matrix based on the same concept is proposed to fit our new problem formulation.

2.1. Brief Introduction of the Fractal Matrix

The fractal matrix was first proposed by Gao et al. [36]. It is a special type of reference matrix which is applied to guide the production of shares in their secret sharing scheme. A fractal matrix is constituted by a lot of fundamental three-dimensional structures called fractal models. Two types of fractal models sized $2 \times 2 \times 2$ and $3 \times 3 \times 3$ are shown in Figs. 2(a) and 2(b), respectively. In each model, the colored boxes represent the fractal elements in the model. There are four and nine elements in Figs. 2(a) and 2(b), respectively. Following the spatial self-similarity property of fractals, a fractal

group can be constructed using fractal models as basic elements. As shown in Fig. 3, the fractal matrix proposed in [36] contains a long series of concatenated fractal groups on its main diagonal.

The elements in a fractal model are so ingeniously arranged that the projections of all elements in a fractal group onto each axial plane constitute a perfect square matrix. In other words, each element in the fractal group is unique and indispensable. This property was exploited to design the (2, 3) SIS scheme in [36]. To achieve a similar purpose, the fractal matrix proposed in this paper is a modified version of the above matrix and will be described in the next subsection.

2.2. The Newly Proposed Fractal Matrix

As discussed in the previous subsection, two fractal models, see Fig. 2, are

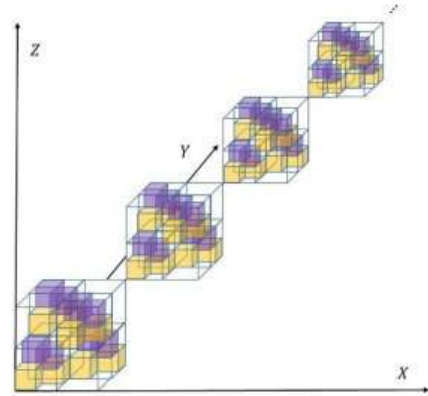


Fig.3: Schematic diagram of the fractal matrix.

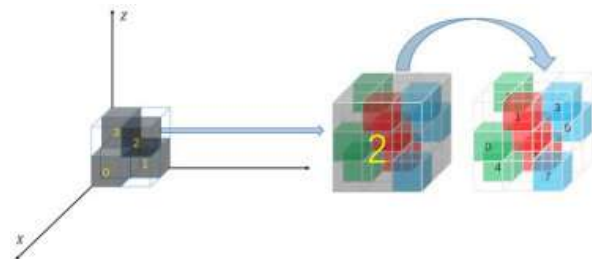


Fig.4: Illustration of the hybrid fractal group

applicable to construct a fractal group. Two factors to be considered are the embedding capacity and the image distortion. A large sized fractal group has a high embedding capacity but the distortion of the shares is also large.

In this paper, we adopt a fractal group of hybrid type, where fractal models of FI sized $3 \times 3 \times 3$ are nested within a fractal model of FO sized $2 \times 2 \times 2$ to constitute a fractal group sized $6 \times 6 \times 6$. As shown in Fig. 4, a fractal group is illustrated in the three-dimensional coordinate system. Assuming we have a three-dimensional matrix $M(x, y, z)$, where $0 \leq x, y, z < 6$, Eqs. (1-4) can be applied to check whether $M(x, y, z)$ is an element in the fractal group. The cubic space sized $6 \times 6 \times 6$ is divided into eight blocks of size $3 \times 3 \times 3$. The coordinates (xO, yO, zO) are numbered block-wise for the outer fractal model, while the coordinates (xI, yI, zI) are numbered pixel-wise for the inner fractal model. Eqs. (1) and (2) determine the outer and inner coordinates. Eq. (3) determines whether $M(xO, yO, zO)$ is an element of the outer fractal model, while Eq. (4) determines whether $M(xI, yI, zI)$ is an element of the inner fractal model. Both conditions should be valid to conclude that $M(x, y, z)$ is an element of the proposed fractal group.

$$xO = \lfloor bx/3 \rfloor; \quad yO = \lfloor by/3 \rfloor; \quad zO = \lfloor bz/3 \rfloor. \quad (1)$$

$$xI = x \bmod 3; \quad yI = y \bmod 3; \quad zI = z \bmod 3. \quad (2)$$

$$yO = (-xO + zO + 1) \bmod 2, \quad (3)$$

$$yI = (-xI + zI + 1) \bmod 3, \quad (4)$$

After determining the configuration of the fractal group, we apply a secret key to initialize the random number generator and start the numbering process. First, we assign a random permutation of 0 to 7 to all elements of each fractal model FI except for the central element, which is reserved for authentication purposes. Next, we assign a random permutation of 0 to 3 to each FO inside the fractal group. For instance, FO (1, 0, 0) 0, FO (0, 1, 0) 1, FO (0, 0, 1) 3, and FO (1, 1, 1) 2 as shown in the

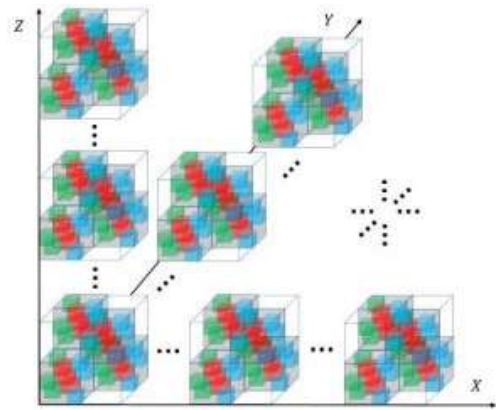


Fig.5: illustration of the newly proposed fractal matrix

figure. So there are a total of 23×22 elements in a fractal group.

The arrangement of fractal groups in the new fractal matrix is quite different from the arrangement in the original version proposed in [36]. The new fractal matrix is illustrated in Fig. 5. The duplicated fractal groups are compactly stacked to constitute a huge cube of $(255)/6 \times (255)/6 \times (255)/6 \times 42 \times 42 \times 74088$ fractal groups in total, which denotes the floor operation. The rest elements outside the range of $[0, 251, 0, 251, 0, 251]$ covered by the fractal groups are left unassigned for authentication.

3. The Proposed (2, 3) Secret Image Sharing Scheme

In this section, we first take an overview of the proposed scheme. Then, the production process of shares, authentication mechanisms, and data extraction process are introduced sequentially.

3.1. Overview

The schematic diagram of the proposed (2, 3)-SIS scheme is illustrated in Fig. 5. The data hider uses three distinct cover images to embed secret images or data and generates three shares, which are visually indistinguishable from their corresponding

cover images. The three shares are then distributed to three participants. To recover the secret data, any two of the participants who combine their shares can completely decrypt the secret. We provide an authentication mechanism to detect tampered shares and prevent cheating events, to be used before the decryption. When all three shares can be obtained, an even stricter authentication mechanism is available.

3.2. Production of the Shares

Assume that the secret data to be shared is converted into a long stream of binary segments denoted by

$S = \{sk \mid k = 1, 2, \dots, n\}$. For each binary segment, $sk_1 \sim 2$ and $sk_3 \sim 5$ denotes the front 2 bits and the rear 3 bits, respectively. Three grayscale cover images with size $H \times W$ are rearranged and denoted as $I^1 = \{p1i \mid i = 1, 2, \dots, H \times W\}$, $I^2 = \{p2i \mid i = 1, 2, \dots, H \times W\}$ and $I^3 = \{p3i \mid i = 1, 2, \dots, H \times W\}$, respectively. The shares to be generated are denoted by I^1 , I^2 , and I^3 , which are the same size as the original cover images and are represented by $I^1 = \{p^1i \mid i = 1, 2, \dots, H \times W\}$, $I^2 = \{p^2i \mid i = 1, 2, \dots, H \times W\}$, and $I^3 = \{p^3i \mid i = 1, 2, \dots, H \times W\}$, respectively.

$\dots, H \times W\}$, $I^2 = \{p^2i \mid i = 1, 2, \dots, H \times W\}$, and $I^3 = \{p^3i \mid i = 1, 2, \dots, H \times W\}$, respectively.

In this share generation phase, firstly, the three-dimensional fractal matrix M is constructed. Then, the pixels in cover images are processed in a prearranged order. Each time, we apply a pixel triplet $(p1i, p2i, p3i)$ as the coordinates and map to the elements in the fractal matrix M $(p1i, p2i, p3i)$. Locate its mother fractal group FG and the precisely mapped fractal element by Eqs. (5-8).

$$Qx = bp1i/6; Qy = bp2i/6; Qz = bp3i/6. \quad (5)$$

$$xt = p1i \bmod 6; yt = p2i \bmod 6; zt = p3i \bmod 6. \quad (6)$$

$$xO = bxt / 3; yO = byt / 3; zO = bzt / 3. \quad (7)$$

$$xl = xt \bmod 3; yl = yt \bmod 3; zl = zt \bmod 3. \quad (8)$$

Its mother fractal group, outer fractal element, and inner fractal element are

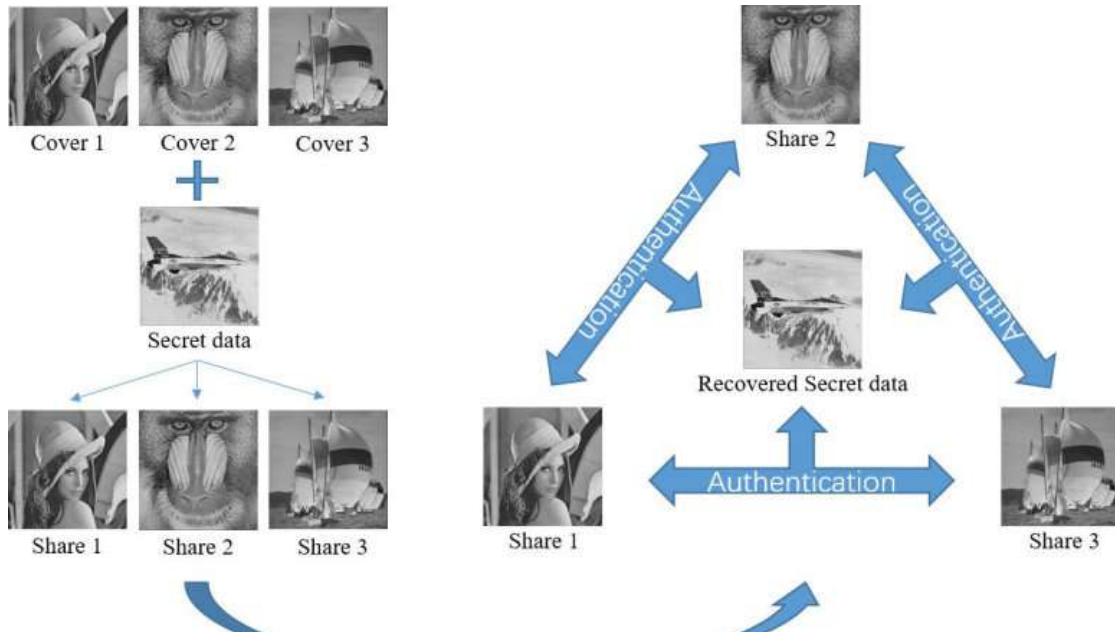


Fig.6: Schematic diagram of the proposed (2, 3)-SIS scheme

denoted by $FG(Qx, Qy, Qz)$, $FO(xO, yO, zO)$, and $FI(xI, yI, zI)$, respectively, where $0 \leq Qx, Qy, Qz < 41$, $0 \leq xO, yO, zO < 2$, $0 \leq xI, yI, zI < 3$. To embed a secret segment sk , we search the mother fractal group $FG(Qx, Qy, Qz)$ to find an outer fractal element $FO(x^{\wedge}O, y^{\wedge}O, z^{\wedge}O)$ and an inner fractal element $F(x^{\wedge}I, y^{\wedge}I, z^{\wedge}I)$ such that $F \subset FO$, $x^{\wedge}O, y^{\wedge}O, z^{\wedge}O = sk_{1 \sim 2}$ and $FI, x^{\wedge}I, y^{\wedge}I, z^{\wedge}I = sk_{3 \sim 5}$. Then, calculate the share pixel triplet $(p^{\wedge}1i, p^{\wedge}2i, p^{\wedge}3i)$ by

$$p^{\wedge}1i = 6 \times Qx + 3 \times x^{\wedge}O + x^{\wedge}I; \quad (9)$$

$$p^{\wedge}2i = 6 \times Qy + 3 \times y^{\wedge}O + y^{\wedge}I; \quad (10)$$

$$p^{\wedge}3i = 6 \times Qz + 3 \times z^{\wedge}O + z^{\wedge}I. \quad (11)$$

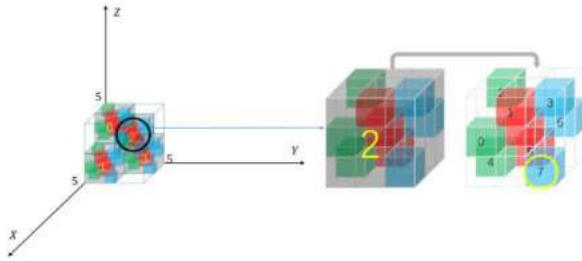


Fig.7: An example of the share triplet determination.

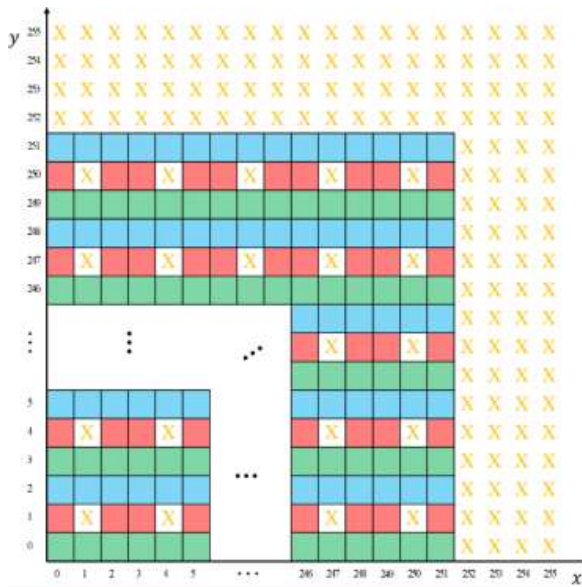


Fig.8: Projection of the fractal matrix on the xy-plane

The share pixel triplet is continually recorded to the shares $I^{\wedge}1, I^{\wedge}2$, and $I^{\wedge}3$ as the embedding proceeds. For example, as shown in Fig. 6, suppose the original pixel triplet is $(0, 2, 5)$ and the secret segment $sk = (10111)_2$. The mother fractal group, outer fractal element, and inner fractal element are $FG(0, 0, 0)$, $FO(0, 0, 1)$, and $FI(0, 2, 2)$, respectively. To embed secret bits $sk_{1 \sim 2} = (10)_2 = 2$ and $sk_{3 \sim 5} = (111)_2 = 7$, we can find $FO(1, 1, 1) = 2$ and $FI(2, 2, 0) = 7$ that satisfy our embedding rule. The share pixel triplet is therefore $(0 + 3 + 2, 0 + 3 + 2, 0 + 3 + 0) = (5, 5, 3)$ as illustrated in Fig. 6.

For the special cases of $Q = p/6] = 42$, i.e., $252 \leq p \leq 255$, the mapped matrix element does not belong to any fractal group. To embed data with minimal distortion, we round it to the nearest index of $Q/41$. The share production algorithm is summarized as follows.

4. Authentication of the Shares

As discussed in Section 3.1, any two of the participants can recover the secret data by sharing their shares. In order to prevent cheating, two authentication mechanisms are provided.

4.1. Authentication with Two Shares

Referring to Section 2.2, the central element of each fractal model and the residual space resulting from fractal group stacking are left unnumbered. The projection of the fractal matrix on the xy-plane is illustrated in Fig. 7. Since the fractal groups are duplications of a fundamental fractal group, the overlapped fractal elements at a particular position corresponding to the same element in the fundamental group. Thus a unique secret binary segment can be determined. The positions marked by 'x' are vacated in the fractal matrix. Therefore, a proper secret embedded share triplet must project to a color position. A share pixel pair $(p1i, p2i)$ which maps to an 'x'-marked position indicates the existence of a tampered pixel. This principle is valid for projections on all axial planes. An authentication mechanism is

Algorithm 1 Production of the Shares

Input:	Cover images I_1, I_2 , and I_3 , where $I_m = \{p_{mi} i = 1, 2, \dots, H \times W\}$, Secret data $S = \{s_k k = 1, 2, \dots, n\}$, Secret key K .
Output:	shares \hat{I}_1, \hat{I}_2 , and \hat{I}_3 , where $\hat{I}_m = \{\hat{p}_{mi} i = 1, 2, \dots, H \times W\}$.
1:	Construct the fractal matrix (see Section 2.2) using secret key K .
2:	For each pixel triplet (p_{1i}, p_{2i}, p_{3i}) ,
3:	Retrieve a data segment s_k .
4:	Calculate Q_x, Q_y , and Q_z by Eq. (5).
5:	For all Q ,
6:	If $Q = 42$, modify to $Q = 41$. End
7:	End
8:	Search in the fractal group $F_G(Q_x, Q_y, Q_z)$ to find $x^*o, \hat{x}^*o, \hat{z}^*o$ and $\hat{I}_1, \hat{I}_2, \hat{I}_3$, such that $F_O(x^*o, \hat{x}^*o, \hat{z}^*o) = s_k^{1 \sim 2}$ and $F_I(\hat{I}_1, \hat{I}_2, \hat{I}_3) = s_k^{3 \sim 5}$.
9:	Calculate the share pixel triplet $(\hat{p}_{1i}, \hat{p}_{2i}, \hat{p}_{3i})$ by Eqs. (9-11).
10:	Record $(\hat{p}_{1i}, \hat{p}_{2i}, \hat{p}_{3i})$ to shares \hat{I}_1, \hat{I}_2 , and \hat{I}_3 .
11:	End
12:	Output shares \hat{I}_1, \hat{I}_2 , and \hat{I}_3 .

devised according to this principle and given as follows.

4.2. Authentication with Three Shares

When we can obtain all three shares, more information can be applied to detect a tampered share. A secret embedded share pixel triplet must be the coordinates of a fractal element. Otherwise, a tampering event is detected. Under such circumstances, we can authenticate the three shares pairwise. When a pair has passed the authentication, the remaining

Algorithm 2 Authentication With Two Shares

Input:	Shares \hat{I}_1, \hat{I}_2 , where $\hat{I}_m = \{\hat{p}_{mi} i = 1, 2, \dots, H \times W\}$.
Output:	Number of tampered pixel pairs N_F .
1:	$N_F = 0$.
2:	For each pixel pair $(\hat{p}_{1i}, \hat{p}_{2i})$,
3:	If $(\hat{p}_{1i}$ or $\hat{p}_{2i} > 255)$, /*check boundary */
4:	$N_F = N_F + 1$.
5:	Else
6:	If $\hat{p}_{1i} \bmod 3 = 1$ and $\hat{p}_{2i} \bmod 3 = 1$, /*check central element*/
7:	$N_F = N_F + 1$.
8:	End
9:	End
10:	End
11:	Print N_F .
12:	If $N_F = 0$, Print "authentication passed." End

Algorithm 3 Authentication With Three Shares

Input:	Shares \hat{I}_1, \hat{I}_2 , and \hat{I}_3 , where $\hat{I}_m = \{\hat{p}_{mi} i = 1, 2, \dots, H \times W\}$.
Output:	Number of tampered pixel triplets N_F .
1:	$N_F = 0$.
2:	For each pixel triplet $(\hat{p}_{1i}, \hat{p}_{2i}, \hat{p}_{3i})$,
3:	If $(\hat{p}_{1i}$ or \hat{p}_{2i} or $\hat{p}_{3i} > 255)$, /*check boundary*/
4:	$N_F = N_F + 1$.
5:	Else
6:	Substitute $(\hat{p}_{1i}, \hat{p}_{2i}, \hat{p}_{3i})$ into Eqs. (5-8) in order to obtain $F_O(xo, yo, zo)$ and $F_I(x_I, y_I, z_I)$.
7:	If $yo = (-xo + zo + 1) \bmod 2$, /*check F_O^* */
8:	If $(y_I \neq (-x_I + z_I + 1) \bmod 3)$ or $(x_I = y_I = z_I = 1)$, /* F_I^* */
9:	$N_F = N_F + 1$.
10:	End
11:	Else
12:	$N_F = N_F + 1$.
13:	End
14:	End
15:	End
16:	Print N_F .
17:	If $N_F = 0$, Print "authentication passed." End

one is a tampered share. When there are no pairs of shares that can pass the authentication, it indicates that at least two shares are tampered with and that the secret is lost. The details of the algorithm are given as follows.

4.3 Extraction of Secret Data

Algorithm 4 Extraction of Secret Data

Input:	Shares \hat{I}_a and \hat{I}_b , where $\hat{I}_m = \{\hat{p}_{mi} i = 1, 2, \dots, H \times W\}$, Secret key K .
Output:	Secret data $S = \{s_k k = 1, 2, \dots, n\}$.
1:	Construct the fractal matrix (see Section 2.2) using secret key K .
2:	Project $F_O(xo, yo, zo)$ and $F_I(x_I, y_I, z_I)$ onto ab -plane to obtain $F_{Oab}(ao, bo)$ and $F_{Iab}(ai, bi)$.
3:	For each pixel pair $(\hat{p}_{ai}, \hat{p}_{bi})$,
4:	Substitute $(\hat{p}_{ai}, \hat{p}_{bi})$ into Eqs. (4)-(6) to obtain $\hat{a}o, \hat{b}o, \hat{a}i, \hat{b}i$.
5:	Extract secret digits by $s_k^{1 \sim 2} = F_{Oab}(\hat{a}o, \hat{b}o)$, $s_k^{3 \sim 5} = F_{Iab}(\hat{a}i, \hat{b}i)$.
6:	Record s_k to S .
7:	End
8:	Output S .

Any two shares that have been passed the authentication can cooperate to extract

secret data. Each corresponding pair retrieved from the two shares can be applied to extract a secret segment. The extraction algorithm is provided as follows.

5. Experimental Results

In this experimental section, we will evaluate the performance of our SIS scheme. The 12 standard grey images of size 512 x 512 used in our experiment are shown in Fig. 8. All programs are implemented with MATLAB R2018a.

5.1. Visual Quality

In order to evaluate the visual quality of the shares generated by our SIS scheme, we use two evaluation indicators, PSNR (Peak Signal-to-Noise Ratio) and SSIM (Structural Similarity).

The formula for PSNR is as follows:

$$PSNR = \frac{255^2 \times H \times W}{\sum_{i,j} (p_i - p_j)^2} \quad (12)$$

where H and W represent the height and width of the two images, and p_i and p_j represent the corresponding pixels in the two images. In general, if the PSNR value of the share and the cover image exceeds 30dB, it is difficult for human eyes to distinguish the difference between the two images.

The formula for SSIM is as follows:

$$SSIM(P, P_r) = \frac{(2\mu_P\mu_{Pr} + c_1)(2\sigma_{PP} + c_2)}{(\mu_P^2 + \mu_{Pr}^2 + c_1)(\sigma_P^2 + \sigma_{Pr}^2 + c_2)} \quad (13)$$

where P and P_r represent the cover image and the share, respectively, μ_P and μ_{Pr} represent the average of P and P_r respectively, σ_P^2 is the variance of P, σ_{Pr}^2 is the variance of P_r and σ_{PP} is the covariance of P and P_r , and c_1 and c_2 are constants used to maintain stability of the formula. c_1 and c_2 can be

obtained by Eqs. (14) and (15), respectively, where L is the dynamic range of pixel values, $k_1 = 0.01$, $k_2 = 0.03$.

$$c_1 = (k_1 L)^2, \quad (14)$$

$$c_2 = (k_2 L)^2. \quad (15)$$



Fig.9: 12 standard grey images

Since the three shares generated by our SIS scheme are only minor modifications to the pixels of the three cover images, the image quality of the shares should be excellent in theory. To confirm this, Fig. 10 shows the first experimental results. The secret image "Airplane" with a size of 256×256 (see Fig. 10(a)) is embedded into three cover images "Coach", "Lena" and "Barbara" with a size of 512×512 (see Fig. 10(b)-(d)). After full embedding, the three high-quality shares are generated with PSNR values of 44.36 dB, 44.39 dB and 44.39 dB, respectively (see Figs. 10(e)-(g)). As shown

in Fig. 10(h), the secret image can be recovered losslessly as long as there are no tampered shares.

The visual quality of the shares generated by our SIS scheme is listed in Table 1, where random bit streams are used as the secret data. In this case, the PSNR value is maintained at about 40.4 dB.



Fig.10: Embedding and restoration results by our SIS scheme.

The embedding capacity is increased to 5 bits per cover pixel pair with only a slight reduction in the quality of the shares. Table 2 compares the features of our scheme with other SIS schemes. SIS schemes based on modern lightweight computing focus on generating meaningful shares with good visual quality and large embedding capability. As shown in the table, the features of reversibility and different cover images are mutually contradictory. Recovering distinct cover

images requires three times the information of a single cover image. It is not worth doing, since the cover images are just used for covering the secret transmission. In addition, the attacker can analyze information from the difference of secret shares generated with a single cover image. Therefore, the single cover image-based schemes are more vulnerable under steganalysis. Our proposed scheme provides three most striking features: the complete difference of the cover images, the extremely high embedding capacity and the (2, 3) secret sharing.

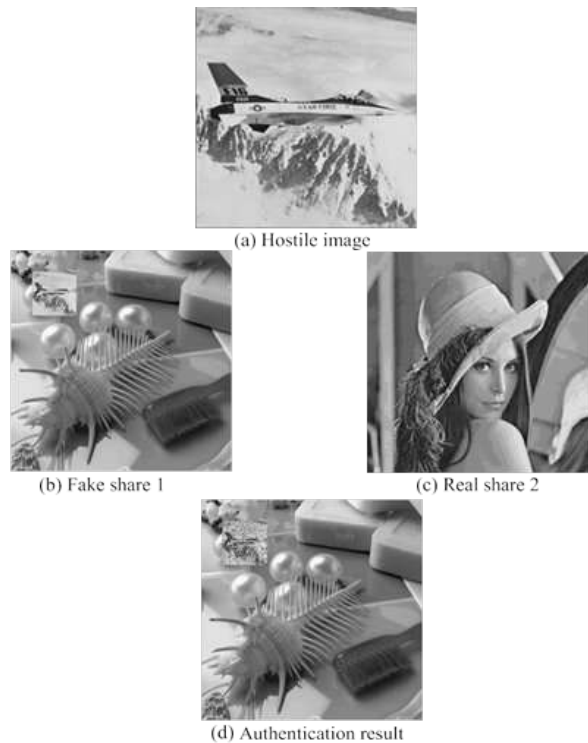


Fig.11: Authentication with two shares.

An advantage of our scheme is that, if one of the shares is invalid, we can use the other two shares to completely recover the secret data, but in the traditional SIS scheme, the secret data cannot be recovered at all. Therefore our scheme provides fault tolerance that the traditional SIS schemes do not have.

Test images (group)	The quality of the shares					
	PSNR (dB)			SSIM		
	Share 1	Share 2	Share 3	Share 1	Share 2	Share 3
Airplane						
Coach	40.37	40.34	40.40	0.95	0.96	0.97
Couple						
Office						
Peppers	40.40	40.43	40.39	0.95	0.96	0.95
Sail						
Baboon						
Wine	40.41	40.40	40.43	0.98	0.95	0.95
Zelda						
Cameraman						
Lena	40.42	40.41	40.42	0.95	0.96	0.97
Boat						

Table 1: Experimental results of the proposed scheme.

5.2. Authentication Ability

In order to test the authentication ability of the two authentication mechanisms proposed in our scheme, we assume that three shares 1, 2 and 3 are distributed to participants 1, 2 and 3, respectively. Suppose the share of Participant 1 is fake, and the shares of Participants 2 and 3 are real. When Participants 1 and 2 want to cooperate to retrieve the secret data, participant 2 can verify whether Participant 1 is dishonest. Here, we use the DR (detection rate) to measure the authentication ability of the two authentication mechanisms. The formula of DR is defined as

$$DR = N_d / N_t, \quad (16)$$



Fig.12: Authentication with three shares.

Features	Gao et al.'s scheme [36]	Chang et al.'s scheme [29]	Li et al.'s scheme [32]	Liu et al.'s scheme [31]	Liu et al.'s scheme [34]	Proposed scheme
Meaningful shares	Yes	Yes	Yes	Yes	Yes	Yes
Reversibility	Yes	-	Yes	Yes	-	-
Different cover images	-	Yes	-	-	Yes	Yes
(k, n) - SIS	(2, 3) - SIS	(2, 2) - SIS	(3, 3) - SIS	(2, 2) - SIS	(2, 2) - SIS	(2, 3) - SIS
Fault tolerance	Yes	-	-	-	-	Yes
Average PSNR	32.96	39.88	49.07	48.72	41.71	40.41
Embedding capacity(bits)	786432	1048576	624215	524288	785525	1310720

Table 2: Comparison with other scheme

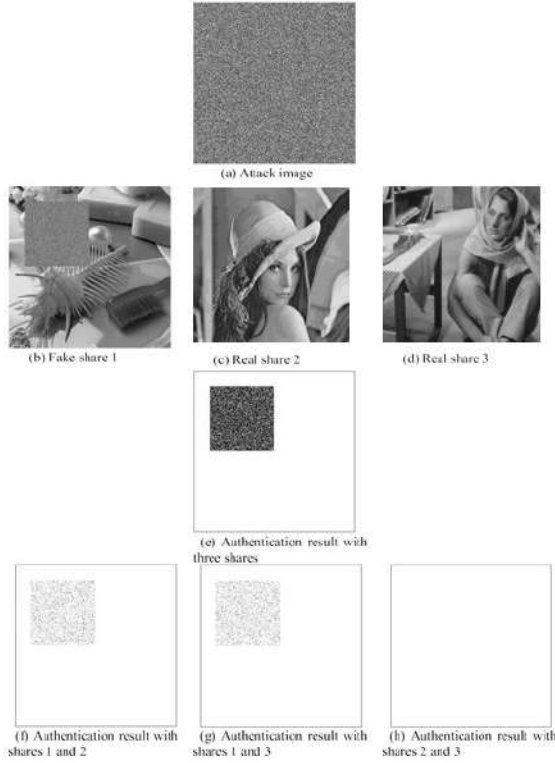


Fig.13: Example of authenticating the three shares pairwise.

where N_d represents the number of illegal

pixels detected by the authentication mechanism, and N_t represents the total number of illegal pixels.

Next, we use two examples to illustrate the two authentication mechanisms which are proposed in our scheme. Referring to Fig. 11, a hostile image "Airplane" (see Fig. 11(a)) is inserted into share 1 in "Coach" (see Fig. 11(b)). The result of using real share 2 "Lena" (see Fig. 10(c)) to verify share 1, the result is shown in Fig. 10(d). The black pixels represent the ones where tampering has been detected.

The second example is shown in Fig. 12, where the hostile image and the tampered share are shown in Figs. 12(a) and (b), respectively. With the help of two real shares (see Figs. 12(c) and (d)), the detection result is shown in Fig. 12(e). From the results, we can see that our scheme has a good ability to detect tampering.

When we can obtain all three shares, more information can be applied to detect a tampered share. A secret embedded share pixel triplet must be the coordinates of a fractal element. Otherwise, a tampering event is detected. Under such circumstances, we can authenticate the three shares pairwise. As shown in Fig. 12,

Test images (group)	DR values					
	Authentication with two shares			Authentication with three shares		
	Standard image	Uniform noise	Special image	Standard image	Uniform Noise	Special image
Airplane						
Coach	0.084	0.097	0.684	0.833	0.836	0.942
Couple						
Office						
Peppers	0.084	0.097	0.685	0.834	0.835	0.942
Sail						
Baboon						
Wine	0.084	0.098	0.684	0.833	0.836	0.943
Zelda						
Cameraman						
Lena	0.084	0.097	0.684	0.833	0.836	0.942
Boat						

Table 3: DR values for different experimental image sets.

when a pair has passed the authentication, the remaining one is a tampered share. When there are no pairs of shares that can pass the authentication, it indicates that at least two shares are tampered with and that the secret is lost.

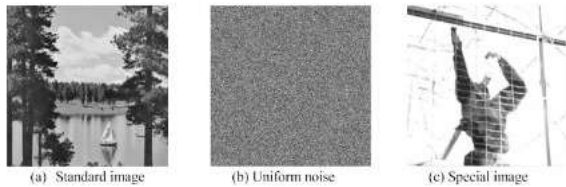


Fig.14: Three images of the different experimental image sets

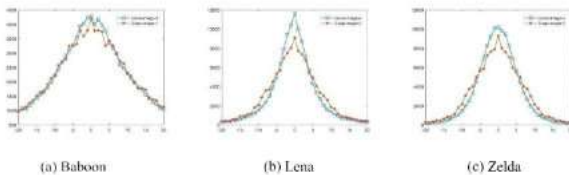


Fig.15: The PDH analysis of the three stego-images and their corresponding cover images.

To investigate the authentication ability of the proposed mechanisms, the DR values for different experimental image sets (see Fig. 13) are listed in Table 3. Since there is a trade-off between embedding capacity and authentication ability for reference matrix-based schemes, the detection rate for authentication with two shares is relatively low. However, the authentication is processed pixel-wise. A small tampering patch consists of hundreds or even thousands of pixels. It is almost impossible to pass the authentication. When all three shares can be obtained, the detection rate can be greatly improved. In fact, the third pixel value of a triplet can be uniquely determined within the whole space of a

particular fractal group.

5.3. Security Analysis

In order to analyze the security of the proposed scheme, pixel-value differencing steganalysis [37] is applied. The histograms of pixel-value difference are shown in Figs. 14, where the test cover images are 'Baboon', 'Lena', and 'Zelda'. The histograms of the fully embedded stego-images are very close to the histograms of the cover images, which indicates that the proposed scheme is secure under the steganalysis of pixel-value differencing.

6. Conclusion

In this paper, a novel (2, 3) SIS scheme based on three distinct cover images and the fractal matrix has been proposed. Our scheme has good security and authentication ability. The scheme has the following features: (1) the three shares are generated based on the three distinct cover images, (2) the shares generated by our scheme have high visual quality, (3) it provides two effective authentication mechanisms, and (4) any two of the three shares can cooperate to extract the secret data losslessly. Experimental results show that the proposed (2, 3) SIS scheme can achieve good embedding capacity with a satisfactory visual quality of the shares.

Our future work will try to generalize the (2, 3) SIS scheme into a (k, n) SIS scheme of arbitrary selected parameter sets, where $k < n$. The principle is to choose a set of elements within a properly sized n -dimensional matrix in a way that their projections into any k -dimensional subspace are not overlapped. This set of elements can be applied as a group for embedding, which guarantees extractability using arbitrarily selected k shares.

7. References

W. Bender, D. Gruhl, N. Morimoto, and A. Lu, "Techniques for data hiding," *IBM Syst. J.*, vol. 35, nos. 3-4, pp. 313-336, 1996.

Z. Ni, Y.-Q. Shi, N. Ansari, and W. Su, "Reversible data hiding," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 16, no. 3, pp. 354-362, Mar. 2006.

C.-C. Chang, "Adversarial learning for invertible steganography," *IEEE Access*, vol. 8, pp. 198425-198435, Oct. 2020.

M. He, Y. Liu, C.-C. Chang, and M. He, "A mini-SuDoku matrix-based data embedding scheme with high payload," *IEEE Access*, vol. 7, pp. 141414–141425, 2019.

J.-H. Horng, S. Xu, C.-C. Chang, and C.-C. Chang, "An efficient data-hiding scheme based on multidimensional mini-SuDoku," *Sensors*, vol. 20, no. 9, p. 2739, May 2020.

C.-C. Chang, C.-T. Li, and Y.-Q. Shi, "Privacy-aware reversible water-marking in cloud computing environments," *IEEE Access*, vol. 6, pp. 70720–70733, 2018.

X. Gao, Z. Pan, E. Gao, and G. Fan, "Reversible data hiding for high dynamic range images using two-dimensional prediction-error histogram of the second time prediction," *Signal Process.*, vol. 173, Aug. 2020, Art. no. 107579.

C.-C. Chang, C.-T. Li, and K. Chen, "Privacy-preserving reversible information hiding based on arithmetic of quadratic residues," *IEEE Access*, vol. 7, pp. 54117–54132, 2019.

C.-C. Chang and Y. Liu, "Fast turtle shell-based data embedding mechanisms with good visual quality," *J. Real-Time Image Process.*, vol. 16, no. 3, pp. 589–599, Jun. 2019.

S. Weng, Y. Shi, W. Hong, and Y. Yao, "Dynamic improved pixel value ordering reversible data hiding," *Inf. Sci.*, vol. 489, pp. 136–154, Jul. 2019.

C.-C. Chang, T.-C. Lu, G. Horng, Y.-H. Huang, and Y.-M. Hsu, "A high payload data embedding scheme stego-images with reversibility," in *Proc. 9th Int. Conf. Inf. Commun. Signal Process.*, Tainan, Taiwan, Dec. 2013, pp. 1–5. C. C. Chang, Y. Liu, and T. S. Nguyen, "A novel turtle shell based scheme for data hiding," in *Proc.*

10th Int. Conf. Intell. Inf. Hiding Multimedia Signal Process., Kita Kyushu, Japan, Aug. 2014, pp. 89–93.

F. Li, Q. Mao, and C.-C. Chang, "A reversible data hiding scheme based on IWT and the SuDoku method," *Int. J. Netw. Secur.*, vol. 18, no. 3, pp. 410–419, 2014.

H. Zhang and L. T. Hu, "A data hiding scheme based on multidirectional line encoding and integer wavelet transform," *Signal Process., Image Commun.*, vol. 78, pp. 331–344, Oct. 2019.

C.-Y. Yang and W.-C. Hu, "Reversible data hiding in the spatial and frequency domains," *Int. J. Image Process.*, vol. 3, no. 6, pp. 373–384, 2010.

X.-Z. Xie, C.-C. Chang, and K. Chen, "A high-embedding efficiency RDH in encrypted image combining MSB prediction and matrix encoding for non-volatile memory-based cloud service," *IEEE Access*, vol. 8, pp. 52028–52040, 2020.

M. Long, Y. Zhao, X. Zhang, and F. Peng, "A separable reversible data hiding scheme for encrypted images based on tromino scrambling and adaptive pixel value ordering," *Signal Process.*, vol. 176, Nov. 2020, Art. no. 107703.

Z. Qian, X. Zhang, and S. Wang, "Reversible data hiding in encrypted JPEG bitstream," *IEEE Trans. Multimedia*, vol. 16, no. 5, pp. 1486–1491, Aug. 2014.

Z. Qian, H. Xu, X. Luo, and X. Zhang, "New framework of reversible data hiding in encrypted JPEG bitstreams," *IEEE Trans. Circuits Syst. Video Technol.*, vol. x29, no. 2, pp. 351–362, Feb. 2019.

F. Huang, X. Qu, H. Kim, and J. Huang, "Reversible data hiding in JPEG images," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 26, no. 9, pp. 1610–1621, Sep. 2016.

Physics Equation of Motion Problem Solver

Nithya P.N., Remya K.R.

Department of Computer Science, Vimala College (Autonomous) Thrissur, Kerala, India | nithya.nimi@gmail.com

Asst. Professor on Contract, Department of Computer, Vimala College (Autonomous) Thrissur, Kerala | remzviji@gmail.com

Abstract

Physics is involved in the day-to-day life of every human being. It helps in understanding how things work from the first principles. Mechanics is one of the major branches of physics that deals with the motion of an object. The concepts in motion come with various word problems which are needed to identify variables and apply them in respective equations. We build a machine solver for physics word problems on the equation of motion. Each problem consists of a question part, describing and setting a formulation part. Here, we describe an NLP-based approach with which a machine can be trained to identify the problem, classify the problem based on 3 equations of motion, apply the suitable equation to the problem, and solve the problem step by step. Python NLTK library is used for entity recognition and the problem is solved using python. A GUI is provided with a text box and the problem is submitted. The model we developed will solve the problem and the solution to the problem is displayed along with all the steps.

Keywords: *Natural Language Processing (NLP), Natural Language Toolkit (NLTK), Bayesian Neural Networks (BNNs)*

1. Introduction

Physics is a branch of Natural science. We can see physics deals with nature and the behavior of matter. Physics deals with solving the problems which arise scientifically. Physics word problems are easy to solve but tricky to understand and identify the variables. Sometimes students find it difficult to identify the variables stated in the word problems and use them in appropriate equations. Physics word problems are taught at secondary and senior secondary levels to develop problem-solving skills.

Humans are no match to computers in terms of computational speed and accuracy. When it comes to problem-solving in various domains, we can see the efficiency in solving problems by the machines to a greater extent provided the problems are in a prescribed syntax of a programming language. This computational ability is seen to be decreasing when we move to human-understandable natural language. Machine-solving word problems is an area of interest for scientists for more than a decade. The help of machine learning and the development of deep learning has made the task of Natural Language

processing an easier task. Natural language processing is involved in various aspects to solve the problems related to daily life

Physics word problem solving is a challenging task because of the semantic gap between the expression and language logic. Our machine solver leverages the commonly followed patterns and positioning of information. It also identifies the information present in the word problem and uses them to solve the problem by substituting it in the relevant Equation of motion. There are 3 equations of motion in Physics which are used to solve the problems on the motion of an object.

2. Literature Survey

The main objective of this system is to solve the physics word problem and provide an accurate solution to the problem. This system helps the user better understand the problem.

An attempt was made by Sowmya S Sundaram along with Deepak Khemani to solve simple word problems using natural language processing. The system takes a word problem described in natural language, extracts information required for representation, orders the facts presented, applies procedures and derives the answer. To elaborate, the problem-solving process begins by processing the question that is expressed in natural language. After processing the question, the information required by the knowledge representation is extracted. Hence, the domain knowledge provided by the underlying representation can also help clear ambiguities faced by the natural language processor.

During the knowledge representation phase, all the sentences are taken into consideration and the events are ordered according to time. Then, the problem is solved using the procedures stored in it. Knowledge Representation Schemas A popular representation idea is to use schemas. How many apples does he have

now?" will trigger the "Transfer in Ownership" schema because of the keyword "forfeit". In the second sentence, the word "forfeit" is such a word that maps to the "transfer-in-ownership" schema. Now the problem is re-examined searching for sentences with "John" and "apple" to match. ROBUST has a disadvantage that, Sometimes, due to the lack of some implicit common-sense knowledge, the problem cannot be solved. the average accuracy of the system is found to be 88.64 percentage [1]

In a similar project called Machine Solver for physics word problems, developed by Megan Leszczynski and Joes Moreira. A machine solver was developed for problems related to free-falling objects under constant acceleration. Each problem has a formulation part, describing the setting, and a question part asking for the value of an unknown quantity. The machine solver consists of two long short-term memory recurrent neural networks and a numerical integrator. The first neural network labels each word of the problem and the second identifies what is being asked in the problem. Here the problems in physics related to the classical mechanics of a point particle in free fall are considered. This domain enables the formulation of a single dynamical system to which all domain-specific physics problems can be mapped. The dynamical system is used to describe the state of the particle, defined by its velocity, changes in the course of time. The classifier in the model is resilient to errors made by the labeler, which performs better in identifying the physics parameters than the question.

The researchers have used Tensor Flow to develop the neural network model for both Labeler and the classifier. The labeler is an LSTM network made of one hidden layer of 10 units. The words act as input into the labeler through an embedding that is randomly initialized and trained simultaneously with the weights and biases which are initially set to zero. The performance of the labeler is assessed through 3 methods: label accuracy, question accuracy, and overall accuracy.

the labeler is trained with TensorFlow's Adam optimizer with initial learning rate of 0.1 and a minimum batch size of 100-word problems was used. The data set consists of 7000-word problems for training, 2000-word problems for validation and 1000-word problems for tests and have obtained accurate results. [2]

In a research work called "learning to automatically solve algebra word problems" Nate Kushman and fellow researchers have presented an approach to automatically learning how to solve algebra word problems. The algorithm reasons across sentence boundaries to construct and solve a system of linear equations that can extract variables and numbers from the problem.

The model can define a joint. Log-linear distribution over full systems of equations and the text. The number and text from the problem is fitted into number slots and the unknown slots are aligned to nouns. Supervised learning is made possible in two scenarios. The first scenario assumes access to the solution of numerical problems and the second scenario is presented with a full system of equations for each problem.

The researchers were able to develop an algorithm which can construct a system of equations by aligning the variable and number to the problem text. The learned model is nearly 69 percentage accurate. The proposed work can be further developed to learn compositional models of meaning for generating new equations and different domains including geometry, physics and chemistry. [3]

In a research work on Simple mathematical word problem solving with deep learning by Sizhu and Nicolas, in their research work, the researchers have developed a model which includes four main steps which include preprocessing the data set. Secondly, three different models are built with a bidirectional LSTM, a bidirectional GRU and a transformer model. Thirdly tuning of Hyperparameters and lastly qualitative analysis of generated output was carried out to study the behavior of

the developed model.

In this work, the deep learning method to solve MWPs problem, the output obtained is one or a set of equations which solve a simple algebra mathematical problem in the text. They have collected and pre-processed more than 45000 question-equation pairs that were collected and preprocessed. Two bidirectional RNNs and transformer models and embedding trained for the data set. The model is successful in generating valid output. Future enhancements can be made to boost the transformation. The model has room for improvement and more feature engineering techniques can be used to obtain the information that is explicitly stated in the problem text. [4]

In another implementation of a similar system, Sonal Srivastava along with other researchers proposed a model for solving physics problems of Class IX NCERT. In the system, the problems are a pre-stored bag of words using training data and standard units are stored and mapped to the keywords. The equations related to different concepts of physics were identified and been and stored. The researchers have also mapped the specific variable to the keywords that are utilized in the corresponding equation. The system consists of an equation solver module and a topic classification module. The equation solver module converts the infix equation to a computer-friendly postfix equation. The obtained equation is parsed and they are processed according to the specific need of the problem. The known values are extracted and put into the corresponding equation to obtain a further simplified equation. The topic classification is done in three steps such as web scraping, calculating TF-IDF for the collection of data, and SVM classifier to classify the topic of the word problem. Here web scraping is carried out using Python's BeautifulSoup Library which creates a dataset with problems of corresponding topics. The system is said to have 78 percentage accuracy and the system is intolerant of spelling mistakes and ways of representing SI units[5].

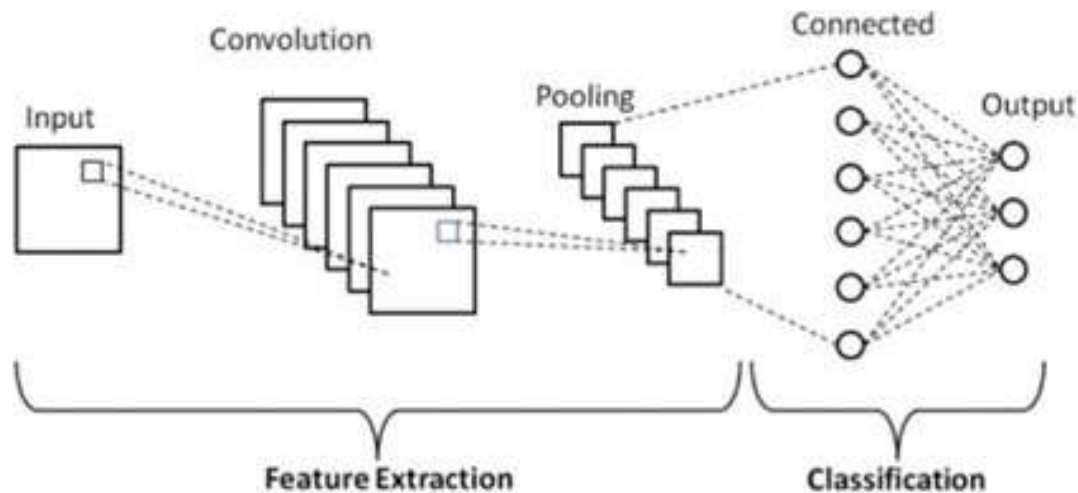


Fig.1: Architecture of Bayesian Convolution neural network

3. Methodology

Solving word problems using Machine learning is always of great interest. A machine cannot understand a word problem as easily as understood by a human. The machine finds it difficult to analyze and process natural language. The development of machine learning, and deep learning technologies have paved the way for NLP which makes it an easy task for the system to process Natural language.

Our system is a neural network with a Bayesian interface which is useful for solving problems in domains where data is scarce. In the Bayesian approach, we use the statistical methodology which can be used to provide a probability distribution attached to it, which is a way to prevent over fitting. The data is fed into the model as a Numpy file, Numpy which is suitable for the word to vector conversion to identify the importance of each word in the given word problem and to provide the words with a specific weight and converted into a categorical value to fit into the model. Pandas is used in the data frame to sort the data according to the weights.

A front end is set for the user interface which accepts the physics problems on the

equation of motion and the required answer is displayed. The front end is developed using the Django framework which is a high-level Python web framework through which a front end can be developed systematically.

4. Model Architecture

We developed a BCNN with 7 input layers and 7 hidden layers. We use Tensor Flow to develop the model. The whole project is divided into 2 parts: Train (1.1) and Test.

4.1. Training phase

During the training phase, the dataset is loaded from google drive. Here we take physics problems that are in the Natural English Language. The data is pre-processed and converted into npz format and then it is loaded into the model. The data flow through train and test is as shown in the figure 1.2.

- Data collection: Data set is prepared by adding problems related to physics. The equation of motion is collected from the NCERT textbook and through online sources are stored in google drive in the form of a csv file. The data stored in the drive is accessed and made accessible to the model for the purpose of training

and testing.

- Separation of Data using NLTK: The data is in the form of a word problem which is written in Natural Language English. Natural Language Processing is made possible by the NLTK library. In this phase we go through the problems and the words are given some weights for extracting the features.
- Store features: The data is in the form of word problems which are written in Natural Language English. Natural Language Processing is made possible by the NLTK library. In this phase we go through the problems and the words are given some weights for extracting the features.
- Train Data: The model is trained on the data set. We have trained the model with 30 epochs.
- Model generation: Once the model is trained, the trained data is stored in a model file which is a h5 file. This model is loaded during test process data.

4.2. Testing Phase

- Load model: Once the model is trained, the trained data is stored into a model file which is a h5 file. This model is loaded during test process data.
- Generate result: The question for which the solution is to be generated is fed into

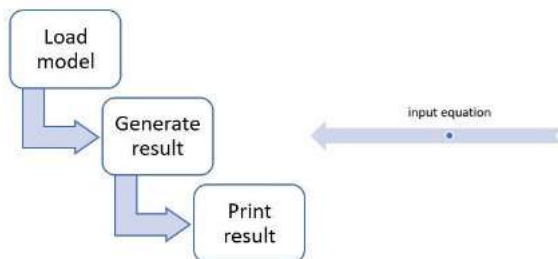


Fig. 1.3: Testing Phase



Fig. 1.2: Training Phase



Fig. 1.4: Testing with problem

the model. Here we have a GUI which consists of a text box. The question is fed into the model and the model identifies the quantity to be calculated in the problem, the respective equation is used to obtain the result.

- Print result: Once the result is obtained it is displayed along with the unit of the particular physical quantity

5. Results and Discussions

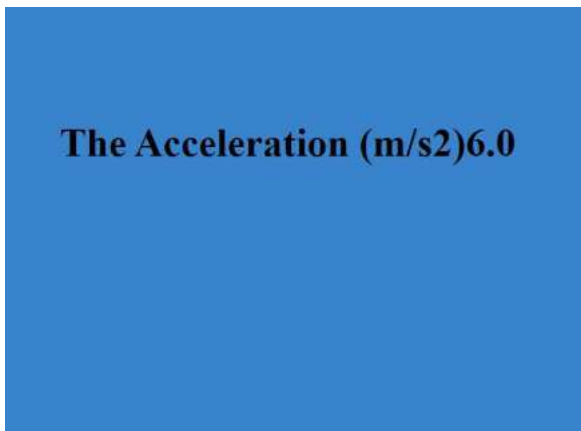


Fig. 1.5: Generating result

We have developed a model for solving physics problems on the Equation of motion. Our model is able to solve

8. References

S.S. Sundaram and D.Khemani, "Natural language processing for solving simple word problems," in *Proceedings of the 12th International Conference on Natural Language Processing*, pp. 394–402, 2015

M. Leszczynski and J. Moreira, "Machine solver for physics word problems," 2016.

N. Kushman, Y. Artzi, L. Zettlemoyer, and R. Barzilay, "Learning to automatically solve algebra word problems," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume1: Long Papers)*, pp. 271–281, 2014.

S. Cheng and N. Chung, "Simple mathematical word problems solving with deep learning," tech. rep., *Technical report, Stanford University*.

problems on objects moving in one dimension, two dimensions or three dimensions. The model is trained with 80 percentage data, 20 percentage for validation, and tested with data from the data set, our model predicts accurate results up to 96 percentage.

6. Conclusion and Future work

The model we developed works well on problems in the equation of motion. We have designed the model with specific language processing formats to identify the keywords and values pertaining to it. The model works well on the problems in the format that is used to train the model, for example, the units of the physical quantities are written in full form such as meter per second for m/s, second for s, kilometer per hour for kmph etc. We hope to extend our technique to handle more general physics problems where units can be mentioned in any format.

7. Acknowledgement

We express our sincere gratitude to God almighty for showering us with all blessings and express our gratitude to all the teaching staff for their valuable guidance and support at each stage of the work. We are also thankful to our parents and family members for the support given in connection with the work.

J. Lampinen and A. Vehtari, "Bayesian approach for neural net-works—review and case studies," *Neural networks*, vol. 14, no. 3, pp. 257–274, 2001.

S. Mandal and S. K. Naskar, "Natural language programming with automatic code generation towards solving addition-subtraction word problems," in *Proceedings of the 14th International Conference on Natural Language Processing (ICON-2017)*, pp. 146–154, 2017.

P. Li, J. Li, and G. Wang, "Application of convolutional neural network in natural language processing," in *2018 15th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP)*, pp. 120–122, IEEE, 2018.

D. G. Bobrow, "Natural language input for a computer problem solving system," 1964.

A. Mukherjee and U. Garain, "A review of methods for automatic understanding of natural language mathematical problems," *Artificial Intelligence Review*, vol. 29, no. 2, pp. 93–122, 2008.

M. Buchanan, *The power of machine learning*. PhD thesis, Nature Publishing Group, 2019.

G. E. Oberem, "Transfer of a natural language system for problem solving in physics to other domains," 1994.

W. Wang and J. Gang, "Application of convolutional neural network in natural language processing," in

2018 *International Conference on Information Systems and Computer Aided Education (ICISCAE)*, pp. 64–70, IEEE, 2018.

D. Zhou, "A new training idea of machine learning in nlp," in *Journal of Physics: Conference Series*, vol. 1861, p. 012083, IOP Publishing, 2021.

J. Forcier, P. Bissex, and W. J. Chun, *Python web development with Django*. Addison-Wesley Professional, 2008.

D. Rubio, *Beginning Django*. Springer, 2017



Pneumonia Category Detection Using Deep Learning

Sreekutty K.V., Resija P.R.

Department of Computer Science, Vimala College Thrissur, Kerala |sreekuttykv88@gmail.com

Department of Computer Science, Vimala College Thrissur, Kerala |resijapr1995@gmail.com

Abstract

In this research, Deep Learning is used for the detection and classification of pneumonia. Pneumonia is a lung infection that causes inflammation in one or both air sacs. Bacteria, viruses, and fungi are only some of the creatures that might cause it. Pneumonia is extremely evident in emerging and underdeveloped nations. Early diagnosis of pneumonia is crucial to ensure better treatment and increase survival rates. The dataset used is the chest X-ray images provided by Kaggle by 5 classes via normal, bacteria, viral, fungal and walking. In this study, using chest X-ray images a Convolutional Neural Network is employed to detect pneumonia infection and its category. Data augmentation strategies are utilized to expand the number of pictures in each class. The VGG16 network is used to develop the CNN model. The performance of testing results attained 88 percentage of accuracy.

Keywords: *Deep Learning, Convolutional Neural Network (CNN), Pneumonia detection, VGG 16*

1. Introduction

A severe respiratory infection known as pneumonia affects the lungs. Small air sacs called alveoli in the lungs of a healthy individual fill up with air when they breathe. When a person has pneumonia breathes, the alveoli are filled with pus and fluid. It leads to painful breathing and limits oxygen intake. The infection can be fatal to anyone, but it poses a serious threat to children, adults, and those under the age of 65. It can range from mild to serious and can sometimes lead to death. The most prevalent cause of pneumonia in adults is bacteria.

Mainly there are four different types of pneumonia. Bacterial, viral, fungal and walking pneumonia. Bacterial pneumonia is caused by bacteria that enter the lungs and multiply there. It might develop on its own or as a result of another sickness, such as a cold or the flu. It affects the immune system. The most harmful variety is viral pneumonia, which is brought on by influenza A and B and the respiratory syncytial virus (RSV). When one or more endemic or opportunistic fungi infect the lungs, it results in a fungal infection. Cough, fever, chest pain, moderate chills, headache, and other symptoms of walking pneumonia can still make you unpleasant.

The disease can be learned using Radiography, CT scan, or MRI. The most often used approach to diagnose pneumonia is chest X-ray imaging.

In this research, we have focused on pneumonia categories and detect the particular type of pneumonia from the chest X-ray images. Convolutional Neural Network is used for this process. Convolutional Neural Networks are a sort of feedforward artificial neural network that takes its connectivity pattern from the visual cortex. In comparison to other classification algorithms, ConvNet requires substantially less preprocessing. While basic approaches require hand-engineering of filters, ConvNets can learn these filters/characteristics with enough training. Through the use of suitable filters, a CNN may properly capture the spatial and temporal dependencies in a picture. CNN has different architectures VGG-16, LeNet, Alexnet, VGG19, and Google LeNet and here we used the VGG-16 model. VGG is a type of convolutional network that is

is a 16-layer deep convolutional neural network.

The following is the key goal of this paper:

Detection and classification of pneumonia through deep learning and convolutional neural networks.

Print the actual type, predicted type and confidence of the corresponding input chest X-ray image.

In Fig. 1, the person has no pneumonia because the predicted label of the chest X-ray in Fig. 1 is normal. Normal stands for the person who has no pneumonia and the type stands for each category.

2. Related Work

Radiography, CT-scan, or MRI can be used to learn about the disease and radiography, or chest X-ray is the most common procedure in a health assessment. Here we use Chest X-ray images to diagnose and treat pneumonia disease. It reduces the number of deaths caused by this disease. A branch of machine learning known as "deep learning" combines representation learning and artificial neural networks. A complex system of algorithms that are modeled on the human brain forms the foundation of deep learning. Consequently, it is possible to process unstructured data, including text, images, and documents.

T. Rajasenbagam, S. Jeyanthi and J. Arun Pandian proposed a method [1] for the detection of pneumonia infection from chest X-ray images using deep convolutional neural network and content based image retrieval techniques. Here the VGG 19 network was used to develop the proposed deep CNN model. In the unseen chest X-ray images, the proposed deep CNN model has a classification accuracy of 99.34 percentage. CNN's performance was compared to that of state-of-the-art transfer learning techniques including AlexNet, VGG16Net, and inceptionNet. The comparison findings reveal that the

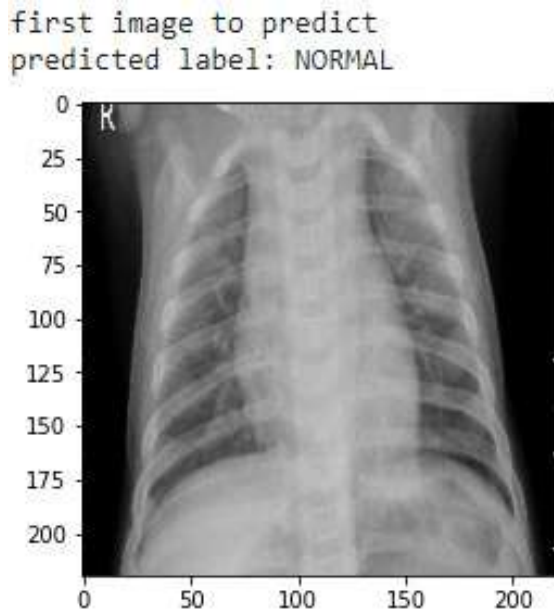


Fig.1: Example

used to classify and locate objects. VGG-16

suggested Deep CNN model outperformed the other strategies in classification.

A. A. Saraiva, D. B. S. Santos, Nator Junior C. Costa, Jose Vigno M. Sousa, N. M. Fonseca Ferreira, Antonio Valente and Salviano Soares proposed an article [2] that explains the comparison of two neural networks for the identification and categorization of pneumonia: the multilayer perceptron and the Neural Network. Cross-validation of k-fold was used to validate the models. The classification models worked well, with the Multilayer Perceptron achieving an average accuracy of 92.16 percentage and the Convolution Neural Network achieving 94.40 percentage.

Another work [3] introduces a study that applies flexible and economical deep learning methodologies, employing six CNN models in predicting and recognising a patient unaffected and affected by the condition. It uses chest X-ray images for the detection of the disease. The six models of CNN here used are GoogLeNet, LeNet, VGG-16, AlexNet, StridedNet, and ResNet-50. Adam is also used as an optimizer in the study, with an adjusted $1e-4$ learning rate and a 500 epoch applied to all of the models. During the training of models, both GoogLeNet and LeNet achieved a 98 percentage accuracy rate, VGGNet-16 achieved a 97 percentage accuracy rate, AlexNet and StridedNet models achieved a 96 percentage accuracy rate, and the ResNet-50 model achieved an 80 percentage accuracy rate. GoogleNet and LeNet models achieved the highest accuracy rate for performance training.

The work of Luka Racic [4] describes the use of machine learning algorithms to interpret chest X-ray pictures in order to aid in the decision-making process in determining the proper diagnosis of pneumonia. It classifies the chest X-ray images into two groups: normal and pneumonia. The study focuses on developing a processing model using a deep learning method based on a convolutional neural network. It achieved a 90 percentage accuracy rate.

Authors in [5] provided a method that examines and compares the identification of lung disease using several computer-assisted methodologies and proposes a new model for identifying pneumonia. A variety of image pre-processing approaches are used for converting raw X-ray images into common formats for analysis and detection. After image processing techniques it uses lung segmentation to acquire the area of interest. CNN, RESNET, CheXNet, DENSENET, ANN and KNN are the different machine learning techniques used here. To detect the presence and absence of pneumonia it uses VGG16 at the end of the classification algorithm. It achieved an accuracy of 96.2 percentage.

Dimpy Varshni, Kartik Thakral, and Lucky Agarwal proposed a method [6] that evaluates the functionality of pre-trained CNN models used as feature extractors followed by different classifiers for the classification of abnormal and normal chest X-Rays. There are two outputs Yes or No, yes stands for pneumonia is present and No stands for the absence of pneumonia. We use analysis to get the best CNN model for the job. The statistical results show that using pre-trained CNN models and supervised classifier algorithms to analyze chest X-ray pictures, specifically to diagnose Pneumonia, can be highly advantageous.

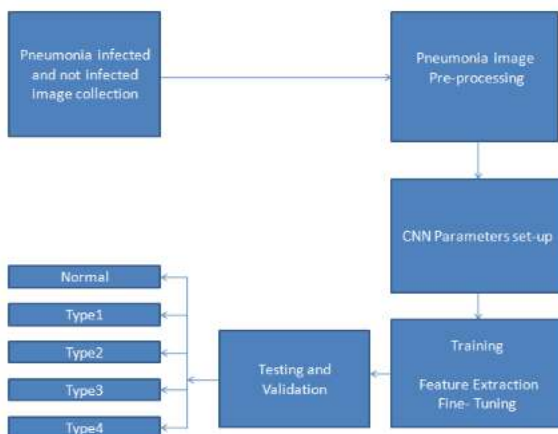
3. Methodology

This paper introduces new models that seek to solve some of the existing research challenges. The existing models are based on binary classification and it only detects if the person has pneumonia or not. The proposed model is based on multilevel classification and it uses deep learning and CNN. Convolution Neural Network is a deep learning algorithm. The Convolution Operation's objective is to take high-level features, such as edges, out of the input image. Traditionally, the first ConvLayer is in charge of capturing Low-Level information such as edges, color, gradient direction, and so on. With further layers, the architecture adjusts to High-Level characteristics as well, giving us a network

that understands all of the photos in the dataset in the same way that humans do. Through the use of suitable filters, a ConvNet may properly capture the spatial and temporal dependencies in a picture. Due to the reduced number of parameters and reusability of weights, the architecture performs superior fitting to the picture dataset. The dataset used for this research is provided by National Clinical Center and is openly available on Kaggle. It contains 1585 images of chest X-rays. Data augmentation is used to increase the number of images in each class. It expands the available dataset by training a deep learning model.

Dataset contains different categories including type1, type2, type3 and type4 and normal. Each type denotes the specific category of pneumonia i.e., bacterial, virus, fungal and walking. Normal indicates that a person has no disease. Lung abnormalities on a chest X-ray will either appear as regions of increased density or as regions of decreased density. Infiltrates and opacification are the attributes used to identify the disease pneumonia from chest X-ray. The normal chest X-ray depicts clear lungs without any areas of abnormal opacification in the image. Bacterial pneumonia exhibits a focal lobar consolidation whereas viral pneumonia manifests with a more diffuse "interstitial" pattern in both lungs.

3.1 System Architecture



The dataset is divided into two, a training set and a testing set. The training ratio is 80 and 20 is the testing ratio. Inside the pneumonia dataset folder, there is training data in the train folder and testing data in the test folder. All images have 220 pixels in height and width. Label encoding is used to change each pixel between 1 to 225. A type of artificial neural network called CNN is used here for image recognition and processing. Higher accuracy can be achieved by enabling shuffling. The dataset loads to CNN. Feature extraction and fine tuning are performed in the training phase. After training the data it generates a model file. In the testing phase, when a new user input arrives it checks which class of the model file has the most similarity with the input and it predicts the output that which category the input image belongs to.

4. Results and Discussion

We have used deep learning and CNN algorithm for the detection of pneumonia. The performance of testing results attained 88 percentage of accuracy. The predicted output contains actual, predicted and confidence. Actual indicates the actual type of pneumonia and predicted denotes the predicted type. It can be normal and has five types. Here we use multilevel image classification and we have 5 classes. These classes are Normal, type1, type2, type3, and type4. Normal stands for the person who has no pneumonia. Types are bacterial, viral, fungal and walking respectively. The following figure shows the predicted output. In this implementation, we fed an x-ray image as input and got the predicted type of pneumonia disease. We can also feed a set of images using a loop as input and detect the categories of pneumonia of a set of X-ray images.

The model is trained using a batch size of 30. We can attain higher accuracy by increasing the number of epochs. A significant hyperparameter for the method is the number of epochs. An epoch in machine learning is one whole run of the training dataset through the algorithm. It specifies, for the entire training dataset,

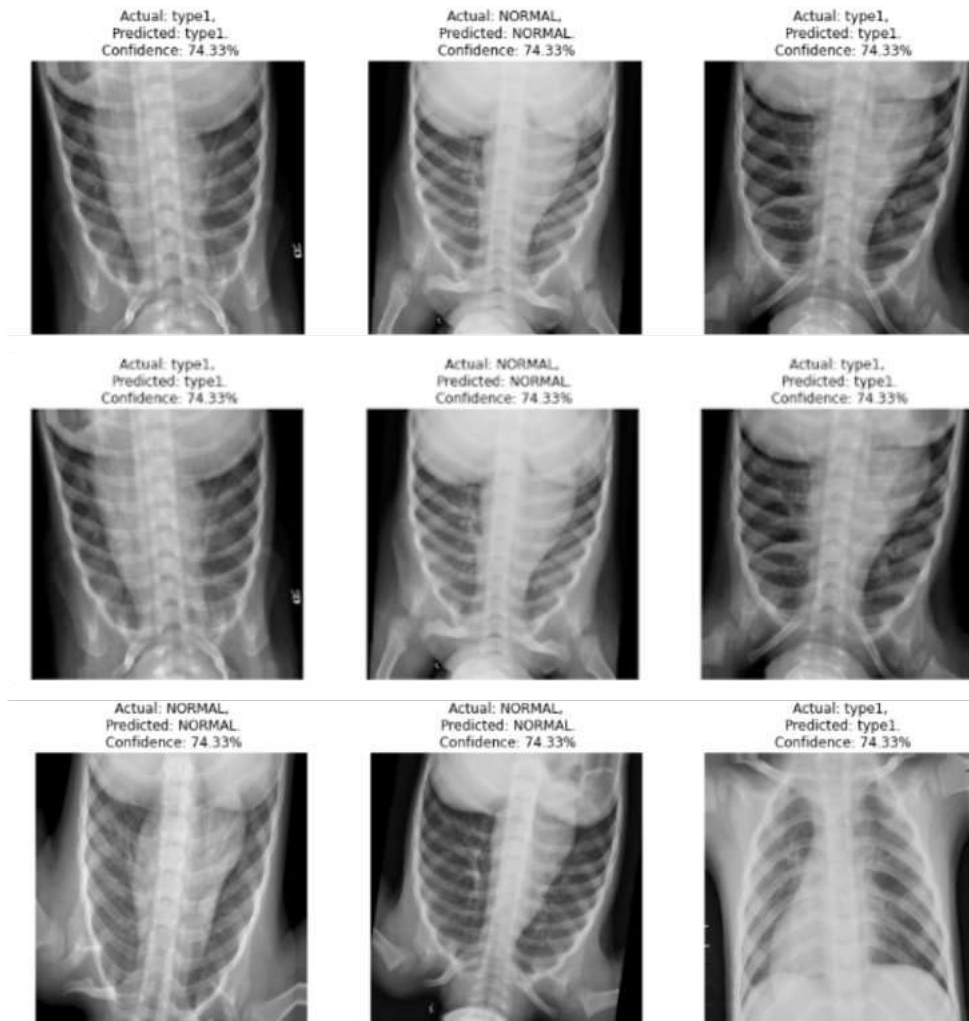


Fig. 2: Output

the number of epochs or full iterations of the algorithm's learning or training phase.

The following figure shows the training and validation losses. It depends on the number of epochs and the losses show small variations when we give the small number for the epochs. If we give a large number of epochs the variations will be high. In the following graph, the number of epochs is 3. Therefore the lines of losses are parallel and there is no variation in the training loss. It is very small. The line of validation loss is parallel to training loss. It is almost consistent throughout the values

and there is a small variation on different stages. The training and validation losses are not parallel to each other. The validation loss is at the peak in some stages and low in others. In Fig. 2.2, the number of epochs given is 60. The lines of training and validation losses show higher variations.

5. Conclusion

Convolution neural network and deep learning play an important role in this paper to find pneumonia. The aim of this research was to detect the disease

pneumonia and its corresponding category from chest X-ray images. It is better to find the category of pneumonia disease more than its normal detection. It

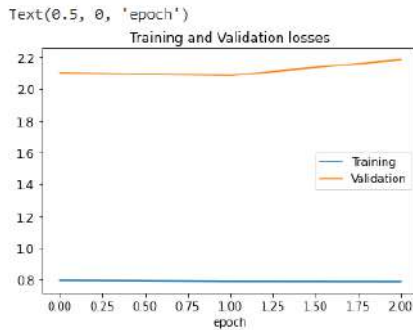


Fig.3: Graph 1

aids in more effective treatment and recovery. Here we used the convolution neural network and its VGG 16 model for this purpose. The proposed method obtains 88 percentage accuracy. It benefits the medical field by providing better treatments to the patients.

7. References

T. Rajasenbagam, S. Jeyanthi, and J. A. Pandian, "Detection of pneumonia infection in lungs from chest x-ray images using deep convolutional neural network and content-based image retrieval techniques," *Journal of Ambient Intelligence and Humanized Computing*, pp. 1–8, 2021.

A. A. Saraiva, D. Santos, N. J. C. Costa, J. V. M. Sousa, N. M. F. Ferreira, A. Valente, and S. Soares, "Models of learning to classify xray images for the detection of pneumonia using neural networks.," in *Bioimaging*, pp. 76–83, 2019.

S. V. Militante, N. V. Dionisio, and B. G. Sibbaluca, "Pneumonia detection through adaptive deep learning models of convolutional neural networks," in *2020 11th IEEE Control and System Graduate Research Colloquium (ICSGRC)*, pp. 88–93, IEEE, 2020.

L. Raćić, T. Popović, S. Sandi, et al., "Pneumonia detection using deep learning based on convolutional neural network," in *2021 25th International Conference on Information Technology (IT)*, pp. 1–4, IEEE, 2021.

A. Tilve, S. Nayak, S. Vernekar, D. Turi, P. R. Shetgaonkar, and S. Aswale, "Pneumonia detection

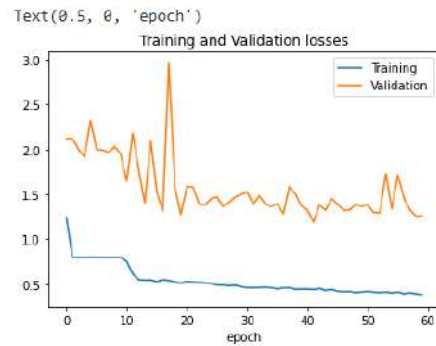


Fig. 4: Graph 2

6. Acknowledgement

We express our sincere gratitude to God Almighty for showering us with all blessings and express our gratitude to all the teaching staff for their valuable guidance and support at each stage of the project. We are also thankful to my parents for the support given in connection with the project.

using deep learning approaches," in *2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE)*, pp. 1–8, IEEE, 2020.

D. Varshni, K. Thakral, L. Agarwal, R. Nijhawan, and A. Mittal, "Pneumonia detection using cnn based feature extraction," in *2019 IEEE international conference on electrical, computer and communication technologies (ICEECT)*, pp. 1–7, IEEE, 2019.

E. Ayan and H. M. Unver, "Diagnosis of pneumonia from chest x-ray images using deep learning," in *2019 Scientific Meeting on Electrical- Electronics & Biomedical Engineering and Computer Science (EBBT)*, pp. 1–5, IEEE, 2019.

S. A. Khoiriyah, A. Basofi, and A. Fariza, "Convolutional neural network for automatic pneumonia detection in chest radiography," in *2020 International Electronics Symposium (IES)*, pp. 476–480, IEEE, 2020.

T. Gabruseva, D. Poplavskiy, and A. Kalinin, "Deep learning for automatic pneumonia detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 350–351, 2020.

V. Chouhan, S. K. Singh, A. Khamparia, D. Gupta, P.

Tiwari, C. Moreira, R. Damaŝeviĉius, and V. H. C. De Albuquerque, "A novel transfer learning based approach for pneumonia detection in chest x-ray images," *Applied Sciences*, vol. 10, no. 2, p. 559, 2020.

T. Rahman, M. E. Chowdhury, A. Khandakar, K. R. Islam, K. F. Islam, Z. B. Mahbub, M. A. Kadir, and S. Kashem, "Transfer learning with deep convolutional neural network (cnn) for pneumonia detection using chest x-ray," *Applied Sciences*, vol. 10, no. 9, p. 3233, 2020.

H. GM, M. K. Gourisaria, S. S. Rautaray, and M. Pandey, "Pneumonia detection using cnn through chest x-ray," *Journal of Engineering Science and Technology (JESTEC)*, vol. 16, no. 1, pp. 861–876, 2021.

N. M. Elshennawy and D. M. Ibrahim, "Deep-pneumonia framework using deep learning models based on chest x-ray images," *Diagnostics*, vol. 10, no. 9, p. 649, 2020.

A. Pant, A. Jain, K. C. Nayak, D. Gandhi, and B. Prasad, "Pneumonia detection: An efficient approach using deep learning," in *2020 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, pp. 1–6, IEEE, 2020.

H. Sharma, J. S. Jain, P. Bansal, and S. Gupta, "Feature extraction and classification of chest x-ray images

using cnn to detect pneumonia," in *2020 10th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, pp. 227–231, IEEE, 2020.

C. J. Saul, D. Y. Urey, and C. D. Taktakoglu, "Early diagnosis of pneumonia with deep learning," *arXiv preprint arXiv:1904.00937*, 2019.

S. R. Islam, S. P. Maity, A. K. Ray, and M. Mandal, "Automatic detection of pneumonia on compressed sensing images using deep learning," in *2019 IEEE Canadian Conference of Electrical and Computer Engineering (CCECE)*, pp. 1–4, IEEE, 2019.

L. Deng, D. Yu, et al., "Deep learning: methods and applications," *Foundations and trends® in signal processing*, vol. 7, no. 3–4, pp. 197–387, 2014.

S. Albawi, T. A. Mohammed, and S. Al-Zawi, "Understanding of a convolutional neural network," in *2017 international conference on engineering and technology (ICET)*, pp. 1–6, IEEE, 2017.

Q. Li, W. Cai, X. Wang, Y. Zhou, D. D. Feng, and M. Chen, "Medical image classification with convolutional neural network," in *2014 13th international conference on control automation robotics & vision (ICARCV)*, pp. 844–848, IEEE, 2014.

Subscription Details

Convergence - International Journal of ICT Academy of Kerala is a peer reviewed International Journal published annually by ICT Academy of Kerala, Technopark, Trivandrum, Kerala.

ANNUAL SUBSCRIPTION FEES	
Non-Members Institutions	Rs. 500/-
Members Institutions	Rs. 300/-
Faculty Members & Students	Rs. 200/-

Ordering Information

If you want to order multiple copies of the Journal, please contact:

The Chief Editor

Convergence, ICT Academy of Kerala

GF-1, Thejaswini, Technopark Campus

T'pura, Kerala - 695 581

Format & Guideline for Publishing the Article - ICTAK Journal

Cover page: The manuscript should be accompanied by a cover page containing the article title, the short title (not more than 5 words and which may be used in all correspondence), the names and affiliations of all the authors (specify order), along with their postal address, phone and fax numbers, and email address. Details of the authors' name and affiliation should not appear elsewhere in the manuscript. In the case of multiple authors, the cover page should indicate the designated corresponding author.

Second Page: The second page should

contain the article title, the short title, the abstract (not more than 100 words), keywords (a maximum of 8 keywords), and an extended summary (not exceeding 300 words).

Body of the article: The recommended length of papers is 8000 - 10000 words, inclusive of tables and figures. Materials may be formatted in Times New Roman, font size 12 and double spaced. All tables and figures are to be serially numbered, sequentially following references to them in the text. All tables and figures are also to be presented in a separate WORD document and file names should clearly specify the paper to which the exhibits belong. All tables and figures should be in black and white only.

ICTAK Journal follows both British & American spelling wherever possible, explanatory theories/concepts and other background material of a historical or collateral nature, and case illustrations/ anecdotal applications should be presented in text boxes to ensure they do not interfere with the flow of the main text.

References

Authors must acknowledge all the sources they have drawn upon, including direct quotations, as well as ideas, concepts, data, and exhibits. Only those references cited in the main text should be listed in the reference list. Source should be stated briefly in the text, following the author-date convention of by the last name and the year of publication, in parentheses. Citations within the text would read, for

e.g. 'According to Pawlak, (1991)...'or '..(Pawlak, 1991)'. These citations should be amplified in a list of references appearing at the end of the paper. The reference list should be in alphabetical and chronological order, and should include complete bibliographical details, as appropriate-the names(s)of the author(s), year of publication, title of the article/ book, name of the journal, details of the publisher, volume and issue number,and individual page numbers, URL of online sources (online journals, magazines, or newspapers) with access date.

The prescribed style of citation is as follows:

Sample book references

Pawlak, Z. (1991). Rough sets: Theoretical aspects of reasoning about data. Norwell, MA:Kluwer Academic Publishers.

Hair, J.F., Black, W.C., Babin, B.J., Anderson, R.E., & Tatham, R.L., (2006). Multivariate data analysis (6th ed.). New Jersey: Pearson Prentice Hall.

Sample reference to chapter in book

Ravallion, M. (2007). Transfers and safety nets in poor countries: Revisiting the trade-offs and policy options. In V. Abhijit, R.B. Banerjee, & D. Mukhejee (Eds.), *Understanding poverty* (pp. 203-230). Oxford University Press.

Copyright & Permissions

Authors must cede copyright of the article

as finally published to ICTAK JOURNAL if it is accepted for publication, and certify that all copyright requirements in respect of material used directly or indirectly in the article have been duly met. Copyright rests with ICTAK JOURNAL in respect of the material submitted for its use and dissemination in any form or medium, individually or in any collection or other configuration, print, audio-video, electronic or otherwise. ICTAK JOURNAL however grants permission to authors for using the submitted material (subsequent to publication in ICTAK JOURNAL) in any printed books or other publications or derivative works authored or co-authored by them. All other usage will be subject to prior written permission of ICTAK JOURNAL.

Further Details

Any correspondence relating to editorial matters and print subscriptions may be addressed to:

The Chief Editor

Convergence, ICT Academy of Kerala
GF-1, Thejaswini, Technopark Campus
T'puram, Kerala - 695 581

Email: editorial@ictkerala.org

Disclaimer: The editors reserve the right to accept or refuse an article for publication, and they are under no obligation to assign reasons for their decision.



Information & Communication Technology Academy of Kerala

(A Govt. of India supported, Govt. of Kerala partnered Social Enterprise)

GF-1 Thejaswini Building, Technopark Campus, Thiruvananthapuram, Kerala 695581

+91 75 940 51437 | 471 270 0811