

# Data Augmentation using Auxiliary Classifier for Improved Detection of Covid-19



Lakshmisetty Ruthvik Raj, Bitra Harsha Vardhan, Mullapudi Raghu Vamsi, Keerthikeshwar Reddy Mamilla, Poorna Chandra Vemula

**Abstract:** COVID-19 is a severe and potentially fatal respiratory infection called coronavirus 2 disease (SARS-Co-2). COVID-19 is easily detectable on an abnormal chest x-ray. Numerous extensive studies have been conducted due to the findings, demonstrating how precise the detection of coronas using X-rays within the chest is. To train a deep learning network, such as a convolutional neural network, a large amount of data is required. Due to the recent end of the pandemic, it is difficult to collect many Covid x-ray images in a short period. The purpose of this study is to demonstrate how X-ray imaging (CXR) is created using the Covid CNN model-based convolutional network. Additionally, we demonstrate that the performance of CNNs and various COVID-19 acquisition algorithms can be used to generate synthetic images from data extensions. Alone, with CNN distribution, an accuracy of 85 percent was achieved. The accuracy has been increased to 95% by adding artificial images generated from data. We anticipate that this approach will expedite the discovery of COVID-19 and result in radiological solid programs. We leverage transfer learning in this paper to reduce time complexity and achieve the highest accuracy.

**Keywords:** CXR, Convolutional Neural Networks, VGGNET, RESNET

## I. INTRODUCTION

Covid-19 is a lung-related disease caused by severe coronavirus respiratory syndrome 2. (SARS-CoV-2). In the final months of 2019, Covid-19 was discovered in Wuhan and has since developed into a global pandemic. Given the absence of a vaccine, therapy, or cure, the only effective way to protect humans from COVID-19 is to contain its spread through swift population tests and isolation of the affected population. To reduce this cavity, a combination of selected health symptoms and a chest X-ray can be used. A radiation chest can be used to indicate the radiologists' covid infection. This enables the development of numerous profound models of learning, and tests have demonstrated that COVID

infections are most likely detected via chest x-ray photos.

Given sufficient data, the revolutionary neural networks have gained the cutting edge of medicine (CNNs). This is completed by the formation of some of the labeled data and the fine tuning of its million parameters. As a result of the large number of CDN parameters that are easily adaptable to smaller datasets, the efficient generalizing process is proportional to the size of the etiquette data. With a fixed number and variation of samples in medical imaging, the most difficult challenge is to obtain fewer data sets. Medical image collection is an expensive and time-consuming process that radiologists and investigators must be involved in. Additionally, the recent COVID pandemic complicates the task of collecting the necessary chest X-ray (CXR) data. With the growth of synthetic data, we propose to alleviate the backdrops. Methods of augmenting data for the purpose of training artificial data are acceptable. In the most part, multidimensionality changes such as image scaling, flipping, conversion, contrast or brightness enhancements, blushing and sharpening, and black and whites balance are used. recent data increase techniques. The traditional data increase is quick, reliable, and accessible. The changes in this increment are sometimes specified as an existing sample is structured into a slightly modified sample. In other words, the increase in conventional data is not invisible. The synthetic increase of data over regular increase limits is a new, advanced form of increase.

## II. METHODOLOGY

Numerous medical imaging techniques make use of the existing system in conjunction with the GAN framework. For the purpose of developing artificial images of pulmonary nodules, the VGG16 multi-dimensional network and the Forward and Backward GAN (F & BGAN) DCGAN-based model have been developed. PGGAN (a growing, growing antiretroviral network) was trained for the purpose of integrating clinical imaging images of retinopathy vascular pathology (ROP) and pre-multimodal glioma images. GAN is used to generate distinct images of the lungs and heart during a chest X-ray examination. The GAN algorithm was used to generate CT images of brain patches connected to a compatible MRI. Additionally, they advise models to utilize the default image. To investigate the relationship between the binary cerebral cortex map and brain MRI images, two GAN networks, Segment and Critic, were developed. GAN training entails establishing the Nash equilibrium of the game. At times, a gradient drop will do this; at other times, it will not.

Manuscript received on August 17, 2021.

Revised Manuscript received on September 26, 2021.

Manuscript published on September 30, 2021.

\* Correspondence Author

**Lakshmisetty Ruthvik Raj\***, Department of Computer Science, Vellore Institute of Technology, Vellore (Tamil Nadu), India.

**Bitra Harsha Vardhan**, Department of Computer Science, Vellore Institute of Technology, Vellore (Tamil Nadu), India.

**Mullapudi Raghu Vamsi**, Department of Computer Science, Vellore Institute of Technology, Vellore (Tamil Nadu), India.

**Keerthikeshwar Reddy Mamilla**, Department of Computer Science, Vellore Institute of Technology, Vellore (Tamil Nadu), India.

**Poorna Chandra Vemula**, Department of Computer Science, Vellore Institute of Technology, Vellore (Tamil Nadu), India.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

GAN training is insecure because we do not yet have a robust recovery algorithm. The disadvantage is that training for these networks is challenging. These networks are used to improve the performance of non-closed-form losses (unlike normal loss functions such as log loss or duplicate error). We proposed that the two categories be denoted by the CNN structure (COVID-CXR and Normal-CXR). Additional data extensions can increase the variety of the database by expanding it. Data augmentation is used to generate chest X-ray images (CXR). Using our proposed CNN and vgg-16 architectures, we integrate CXR artificial images created through data addition.

The following contributions to this study are:

1. Suggest that in the production of artificial CXR images, we apply the concept of data betterment.
2. Establish a COVID-19 acquisition, CNN-based model
3. In order to detect the improved availability of COVID-19, the addition of training data by the CNN model is used. When CNN is trained in actual data and performance improvements, an improvement to the performance division from 85 to 95 percent is recorded.

### Benefits:

The suggested method enables small-scale increases in COVID-19 accuracy via the use of synthetic x-ray images. Alone, with CNN distribution, an accuracy of 85 percent was obtained. The most frequently used imaging modality for diagnosing COVID-19 patients is chest x-ray imaging. We can determine whether or not this individual has the coronavirus by examining CXR pictures.

## III. ALGORITHMS

### A. Classification

Data analyzed which did not have a label — we did not know the sample class to which (known as unregulated reading). The question supervised, on the other hand, works with labeled data, in which they know each sample's various classes. If we predict which class the sample belongs to, we call this a classification problem. SciKit-Learn offers several algorithms, In this section, we will see support vector. A new sample must be dangerous or harmless, according to the characteristics of the new undetectable sample, based on the vector machine support model.

You will see that the SVM model predicts the error of novel invisible samples in the test set very successfully. The classification function can be used to calculate this with multimetric printing accurately. Clarity, Reminder, and school F1 Here, each class is shown as  $F1 = 2 \cdot \text{Enter} \cdot \text{Remember/Clarify} + \text{Reminder}$ . Support column calculation of sample number per class.

A powerful separation tool is the Vector Machines support. They function well in large environments, even if the number of features exceeds the sample number. However, the number of examples of your working time is quadratic, making it difficult to train large data sets. Quadratic means that it will take longer to train if you enlarge the database ten times.

You will finally find that there are 30 components in the breast cancer database. This makes visualization or editing of data difficult. We can use a process known as size reduction to look at the most visible data.

### B. Dimensionality Reduction

The reduction of size is another crucial way of learning by machines and data science in general. We'll also look at the database on breast cancer in Wisconsin in this example. There are more than 500 samples in the database, each with 30 traits. The features are associated with photographs of a good model for breast tissue, and the features explain the features of the photographs. Every feature is of true value. The target variance is a different value (negative or negative) and is therefore classified as data.

You will remember the example of the Iris by setting a data distribution matrix in which each element is built-in for possible interactions on every other element in the database. By exploring this structure, you may discover features that will divide the database into groups. Because the database has only four elements, we were able to organize each feature easily. However, as the number of traits grows, this is gradually declining, especially when

You consider the genetic model that has more than 6000 traits. Principal Component Analysis, or PCA is another way to manage large-scale data. PCA is an uncontrolled database-size reduction algorithm. You may want to reduce your data to 2 or 3 levels for planning purposes, for example, and PCA allows. Do this by combining the fundamental features you

It can use in your data organization. PCA is an unchecked algorithm. You provide your data X, and you indicate the number of items to be reduced to the so-called data sample:

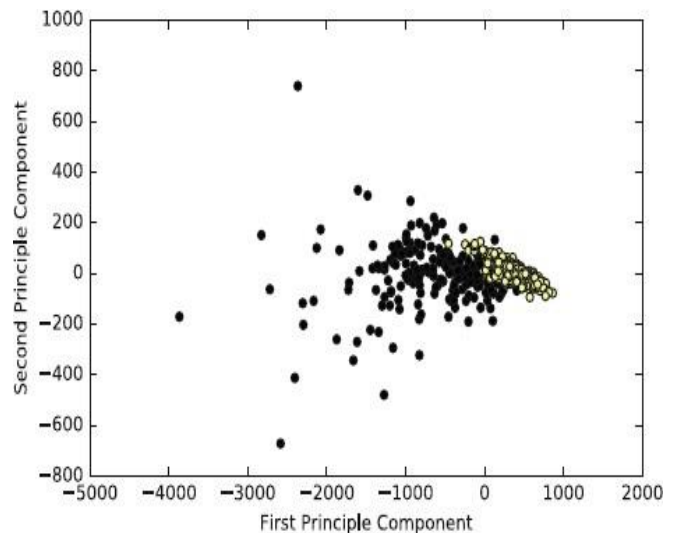


Figure 1: Representing dimensionality reduction

### C. Neural Networks and Deep Learning

Although we will nevertheless discuss an example of how the Kera is one of the most popular deeper learning frameworks, far beyond the chapter defining nerve networks and comprehensive learning. In this section, we will develop a simple neural network to distinguish the previously described Wisconsin breast cancer database. In deep neural networking and learning algorithms, images are often classified — most of them are used in recognizable neural networks.

However, it can also be used for text or table-based information. We build a standard feed, a closely linked network, and separate the cancer database in the text to reflect the use of the frame.

In this example, we also use a database of 30 characteristics and 569 samples for breast cancer in Wisconsin. In order to make the neural network a major challenge, we will use only 50 percent of the total data in the training set. Then we will test the remaining 50% of the neural network.

#### IV. DEEP LEARNING FOR OBJECT DETECTION

##### A. CNN Classifier

In this project, the acquisition of a helmet, scarf, and mask discovery, was done by CNN (Convolutional Neural Networks). The program is trained and tested with images of human hats, scarves, and masks and is used to determine whether a person has a face mask or not.

CNN is a neural network type used for the detection and classification of images. CNN uses supervised learning. CNN contains filters or neurons with biases or weights. Each filter picks up a specific input and changes the received input. The CNN separator has four layers; Convolution, integration, Linear Unit Layer (ReLU), and fully integrated layers.

##### B. VGG16 Model

Transfer learning is usually a process that applies a model trained in one problem to a second related problem in one way or another. During a profound study, transfer learning is the first way to learn a model of a neural network about problem-solving. For the new model trained in this instance of the problem, one or more layers of the pre-trained model shall be applied. Transfer learning has the benefit of cutting the neural network model's training time and can result in a lower overall performance bug. Reusable hand weights can be used as a base for the workout and can be adjusted to solve a new problem. Transfer learning is used as a form of initialization of weight. This can be useful if there is more information on the first related issue, and the similarity in troubleshooting can be helpful in both cases.

##### C. Residual Networks (RESNET)

To lower the failure rate, all of the following machines use multiple layers in a deep neural system after winning the 2012 ImageNet competition from the first CNN architecture (AlexNet). It applies to a small number of layers, but when we increase the number of layers, there is one common problem when reading the so-called vanishing/exploding gradient.

Therefore, the gradient is null or higher. As the number of layers increases, the level of training and test errors increases.

ResNet is launching a new model called Residual Network, proposed by Microsoft Research researchers in 2015.

##### D. Residual Block

This model introduces a concept known as a residual network to address the vanishing/exploding problem in gradients. In this network, we use Skip Connections. Skip connections can be connected directly to the output and overlap with several layers of training. The way behind the network reads the map below, and the network matches the remaining map. The path

behind the network reads. Therefore, let the network equal, instead of  $H(x)$ , to make the first map:  $H(x) - x H(x) = F(x) + x$

If a layer impedes the performance of the model, the benefit of this type of skip connection is the regularization. This causes deep neural network training without the gradient problems caused by the vanishing/exploding. The document was written in 100 to 1000 layers on the CIFAR-10 database. A similar method is used by these networks, called the "highway network," to overcome connections. This skip link uses parametrical gates in a manner similar to LSTM. These doors determine how far the connection passes. However, this structure was no better than the architecture of RESNET.

#### V. DATA SET

Any machine learning project is based on data. The second phase is a complex one involving data collection, selection, processing, and analysis. The steps can be split into further tasks.

##### A. Data collection

The data analyst needs to take the lead and guide the machine in its implementation. The task of a data analyzer is for the collection, interpretation, and analysis of data, using mathematical strategies, ways, and resources.

In this category, the 'more, the better approach.' Some data scientists believe less than one-third of the collected data could help. Estimating which part of the data will give the most accurate results until the model training begins is difficult. This is why all data – internally and externally, organized and uncluttered – is important to collect and store.

Internal data collection tools are based on the infrastructure of industry and business. For example, those who only run an online firm and want to launch a customization campaign can try out a good source of internal data from web analytics tools such as Mixpanel, Hotjar, CrazyEgg, the most popular Google statistics, etc. Keeps user information and online behavior data: time and duration of visits, pages or objects viewed, and location.

Companies can also complete public records for their information. For instance, Kaggle, Github, and AWS provide free information on analytics.

##### B. Pre-processing of data

The purpose of preprocessing data is to transform raw data into a machine learning form. Data scientists can accurately obtain results using systemic and clean data from the used machine-learning model. The process is complemented by formatting, cleaning, and modeling data.

##### C. Formalization of data

When people obtain data from different sources, the importance of data formatting increases. Data scientist's first work on recording formats. The technician examines if each attribute's variables are recorded in the same way. Examples of variables are product and service headings, prices, date formats, and addresses.



The data consistency principle also applies to numerical annotations.

### D. Cleanup of data

This set of procedures enables data inconsistencies to be removed and corrected. A data scientist may use imputation techniques to fill in lost data, e.g., to enter missing values with meanings. The expert also finds the data, which does not match the distributions for the data significantly. If there is an error in the data, the data administrator will remove or, if possible, correct it. The deletion of incomplete and useless data items also includes this category.

### E. Anonymization of data

A data scientist sometimes has to set words or attributes which are sensitive information (e.g., when working with health care and banking data).

### F. Samples of data

More time and computer power are required for larger datasets. If the database is too large, it is possible to use a simple data sample. This method is used by a data scientist to select a small but representative sample of data in order to build and use models extremely quickly while producing accurate results.

## VI. PREPROCESSING OF IMAGE

The processing of pictures is divided into analog and digital picture processing. Processing of digital images Computer algorithms is used in digital images to process images. Digital image processing has more advantages than analog image processing as a small platform for digital signals processing. It can be employed using a vast array of algorithms for input data. Digital imaging (functions) is intended to increase the image data, removing unwanted distortions and/or enhancing other key image functions, to make our AI-Computer vision models more proficient.

### A. Pictures

In this stage, we save a variable from the path to our image database and create folders that include pictures.

### B. Picture Size

In this phase, we create two features that show images to visualize the transition, one showing an image as well as the other showing two images. We then create a processing function that accepts images only as a parameter. The size of some pictures captured and fed into our AI algorithm varies accordingly, so for all images of our AI algorithms, we should define a basic size.

Augmentation of data:

### C. Augmentation of Data

Computer viewing tasks like image classification, object detection, and its subdivision were among the most popular profound learning applications. Data additions can be used effectively to train DL models in such programs. Some of the simple changes applied to the image are; Geometric modifications such as Search, Rotate, Convert, crop, measure, and color space changes such as color input, fluorescent light, and sound injection.



**Figure 2: Data augmentation of various images**

The geometric modification works best when the positional bias trend is present in images such as the data used for facial recognition. Modification of color space can help to deal with challenges connected to light or brightness in images. Data splitting

### D. Data Splitting

The data set used in machine learning projects is divided into three parts. Those parts are training, testing, and validation.

### E. Training Set

A data scientist uses a model training set and sets his/her parameters, which can be learned from the data.

### F. Set of Tests

A test set is required for the testing and standard performance of a professional model. The latter refers to the ability of the model to identify patterns in new data, which cannot be seen in training data. In training and testing, it is important to use different subsets to prevent the model equally since the standard is incapable of performing.

### G. Modeling

At this stage, we train numerous models for getting accurate predictions for our particular data set.

Model training

### H. Model Training

We train models with limited data sets in this step. Fast ai numerous layouts, making transfer learning very easy to use. We can use pre-trained models that work for most applications/data sets to create a convolutional neural network (CNN) model. We will use RESNET building because the data sets with more problems are faster and more accurate. 18 is the number of neural network layers in resnet18. We also transmit measuring metric quality metrics by authentication of the model of the set. We use error rate, which says how often the model predicts the wrong way.

## VII. RESULT AND CONCLUSION

Even if the best COVID-19 test for RT-PCR is considered, it takes time to make a decision because of the serious mistakes in the results. The results are highly positive. As the medical imaging of human lungs, such as X-rays and CT chest scans, is the best alternative according to researchers due to the false-negative results.

Chest X radiation is low-cost and low dose radiation available and easily used in general or local hospitals. Chest X radiation is available. This review presents a detailed study of existing solutions based primarily on DL's COVID-19 Test Result Strategy in the early stages. This study provides a great understanding of scientists' and policymakers' thinking processes - not only during wave times but during vaccinations that might need to be tested in real-time. However, a lack of information is an obligatory challenge in order to achieve effective results in real-time. In this review study, many solutions were presented and discussed to provide further insights into future trends and future diseases that could tackle the problem of data loss. We believe that with more public knowledge about the information like datasets and updated strains, better approaches can be developed to accurately detect and evaluate COVID19.

	model	Test-Accuracy
0	CNN	0.9666
1	VGG-16	1.0000
2	RESNET	1.0000

Figure 3: Performance metrics

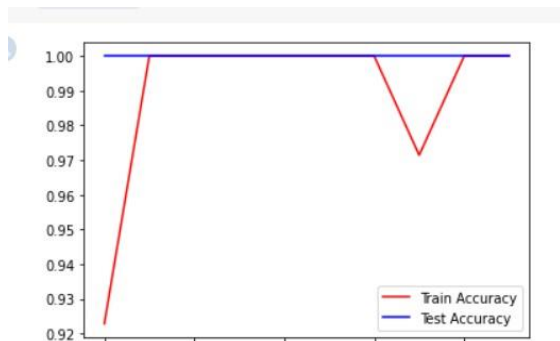


Figure 4: Train Accuracy vs. Test Accuracy

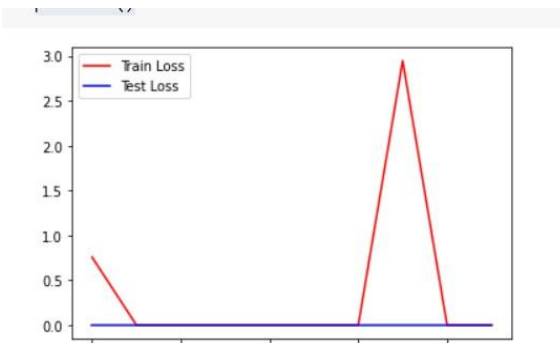


Figure 5: Train Loss vs. Test Loss

REFERENCES

1. S. Wang, J. Sun, I. Mehmood, C. Pan, Y. Chen, and Y. Zhang, "Cerebral micro-bleeding identification based on a nine-layer convolutional neural network with stochastic pooling," *Concurr. Comput. Pract. Exp.*, vol. 32, 2020.
2. S. Wang, C. Tang, J. Sun, and Y. Zhang, "Cerebral Micro-Bleeding Detection Based on Densely Connected Neural Network," *Frontiers in Neuroscience*, vol. 13, 2019.
3. C. Kang, X. Yu, S. Wang, D. S. Guttery, H. M. Pandey, Y. Tian, and Y. Zhang, "A Heuristic Neural Network Structure Relying on Fuzzy Logic for Images Scoring," *IEEE Transactions on Fuzzy Systems*, pp. 1–1, 2020.
4. G. Litjens, "A survey on deep learning in medical image analysis," *Med. Image Anal.*, vol. 42, pp. 60–88, 2017.
5. H. R. Roth, L. Lu, J. Liu, J. Yao, A. Seff, K. Cherry, L. Kim, and R. M. Summers, "Improving Computer-Aided Detection

- Using\_newlineConvolutional Neural Networks and Random View Aggregation," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 1170–1181, 2016. [Online]. Available: 10.1109/tmi.2015.2482920;https://dx.doi.org/10.1109/tmi.2015.2482920
6. H. Greenspan, B. van Ginneken, and R. M. Summers, "Guest Editorial Deep Learning in Medical Imaging: Overview and Future Promise of an Exciting New Technique," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 1153–1159, 2016. [Online]. Available: 10.1109/tmi.2016.2553401;https://dx.doi.org/10.1109/tmi.2016.2553401
7. N. Tajbakhsh, J. Y. Shin, S. R. Gurudu, R. T. Hurst, C. B. Kendall, M. B. Gotway, and J. Liang, "Convolutional Neural Networks for Medical Image Analysis: Full Training or Fine Tuning?" *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 1299–1312, 2016. [Online]. Available: 10.1109/tmi.2016.2535302;https://dx.doi.org/10.1109/tmi.2016.2535302
8. A. Mikołajczyk and M. Grochowski, "Data augmentation for improving deep learning in an image classification problem," 2018 International Interdisciplinary Ph.D. Workshop (IIPhDW), pp. 117–122, 2018.
9. L. Engstrom, D. Tsipras, L. Schmidt, and A. Madry, "A Rotation and a Translation Suffice Fooling CNNs with Simple Transformations," *ArXiv*, vol. abs/1712.02779, 2017.
10. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, . . Bengio, and Y, "Generative adversarial nets," *Advances in neural information processing systems*, pp. 2672–2680, 2014.

AUTHORS PROFILE



**Lakshmisetty Ruthvik Raj** is currently an undergraduate at Vellore Institute of Technology Vellore pursuing a Bachelors's of Technology in Computer Science and engineering. He has worked and developed various projects in the field of Software engineering, Internet of things(IoT), Data Analytics, Computer Vision, and Natural language processing and dealt with a significant number of projects, each having its unique problem. Furthermore, he is optimistic to provide solutions regarding environmental issues, especially in Health care, Agriculture by blending and collaborating his knowledge.



**Bitra Harsha Vardhan** was born in 2000 and is currently pursuing his Bachelors' of Technology in Computer Science and Engineering at Vellore Institute of Technology, Vellore. He completed many subjects in his undergraduate with exceptional grades, such as Natural Language Processing, Image Processing, Data Structures and Algorithms, and many more. He gained experience working on various scenarios with employees in an internship dealing with various user-based systematic computer systems by resolving their problems. Finally, he is looking forward to work on environmental concerns using the latest technological tools along with thoughtful brains.



**Mullapudi Raghu Vamsi**, born in 2000 and currently pursuing my Bachelors in computer science and engineering at Vellore institute of technology, Vellore. He completed his higher secondary education in 2018. His research interests include Artificial Intelligence and Machine learning.



**Keerthikeshwar Reddy Mamilla** was born on 2001 and is currently pursuing his Bachelor of Technology in Computer Science and Engineering at Vellore Institute of Technology, Vellore. He has a good knowledge of research in various fields. He has successfully completed various subjects such as Artificial Intelligence, Database Management Systems, Operating Systems and many more with outstanding grades. He has dealt with various research-based projects. He is a highly motivated individual, who can work on any kind of latest technologies. His constant focus on technologies has driven him to a greater extent. Finally, he is interested to join a group people with same ideology to work on latest technologies.



## Data Augmentation using Auxiliary Classifier for Improved Detection of Covid-19



**Poorna Chandra Vemula** born in 2001 and is currently pursuing his Bachelors' of Technology in computer science and engineering at vellore institute of technology, vellore. He worked on various research projects in the areas of Data Analytics, Machine Learning, NLP and Computer Vision. He has experience working on real world applications dealing with large scale systems, and built numerous apps taking into account current demands of the stakeholders. He also has mentorship experience, explaining concepts in the field of Data Science. He is optimistic about the future developments in AI and looking forward to collaborating on solving problems primarily in the areas of Health, Agriculture, Education and sustainable development using Artificial Intelligence.