# Deep Learning-Based Intelligent Analysis of Medical Big Data

**Reenu Maria Binoy[1], Parvathy Prakash[2] , Dayana Thomas[3] , Devika R Nath[4] ,Gigi Joseph[5]**

*[1,2,3,4] PG - MCA, Mahatma Gandhi University, Kottayam, Kerala*

*[5] Assistant Professor, Kristu Jyoti College of Management and Technology,Kerala*

**ABSTRACT**

Data-driven medical big data analysis methods have arisen as the times have demanded, assisting in the intelligent identification of medical health. With the widespread use of computer technology, medical health data has also expanded considerably. It is still challenging to evaluate medical big data, nevertheless because of the heterogeneous format of the data, the large number of missing records and the amount of noise. While deep learning constructs a hierarchical model by imitating the human brain, traditional machine learning techniques are unable to efficiently harvest the rich information included in medical big data. It boasts strong automated feature extraction, intricate model building and effective feature expression among other significant features.

It is a deep learning technique that takes features from the original medical imaging data at every level, starting at the lowest. For the purpose of intelligent illness identification and diagnosis, this research builds a deep learning-based data analysis model for medical pictures and transcripts. The model chooses and optimises model parameters using a vast amount of medical big data. It also automatically picks up on the pathological analysis process used by doctors or medical researchers. Finally, the model conducts disease judgement and effective decision-making based on the analysis findings of the medical big data. The outcome of the trial demonstrate that the approach can evaluate large amounts of medical data and achieve early illness detection. In addition, it may evaluate the patient's physical health state based on their physical examination data and forecast their future risk of contracting a certain disease. Greatly lower the workload for physicians or researchers in the medical field and increase the effectiveness of their effort.

**Keywords - Medical Big Data, Analysis, Deep Learning, Intelligent Recognition.**

## 1. INTRODUCTION

The use of medical big data recording systems in hospitals has grown steadily more widespread with the expansion of the economy, computer technology and medical information processing technology. The most crucial component of a medical big data study is medical health records. In order to prevent and control a specific disease before it manifests, properly anticipate how the disease will evolve, and identify more high-risk individuals, clinicians can study the condition in accordance with the available medical big data.

Firstly, photographs based on the patient's examination indicators must be recorded and collected for use in the study of medical big data. Then, for illness analysis and early diagnosis, the difference between the normal image and the patient's image is examined. Numerous conventional medical large data analysis techniques, such as logistic regression, random forest algorithm, support vector machine, and others, have been researched. Deep learning's advancement has made it increasingly successful to use this technique for data mining, computer vision, and natural language processing. Additionally, using the task's preferred way to extract characteristics has grown more challenging. It effectively captures long-range relationships in the data and extracts deep and abstract characteristics from it, enabling fast analysis of both picture and text data. Deep learning approaches perform better in big data analysis than conventional methods with the growing number of medical pictures and transcript data. Additionally, they demand fewer time and computer resources for feature extraction and data preparation. The current deep learning model cannot examine the connections between the various sections of the picture and the subtleties therein, nor can it handle transcripts of arbitrary sizes, whether the analysis is performed on medical image data or transcript data. The current deep learning model is not as robust. Additionally, the model's architecture must be altered as the disease's

features change, which requires a lot of time and labour and restricts the model's flexibility. Additionally, it is challenging to use in clinical settings.

In particular, the following can be said about our paper's technical contributions:

First: The challenges, state of the art of study, and research methodologies of deep learning in medical data identification are elaborated and compiled in this work. The relevant deep learning models and methodologies are investigated. Additionally, the network's design offers a reference point.

Second: For the most typical medical pictures and text records in medical big data, various neural network topologies have been developed. To achieve intelligent illness detection and prediction, data analysis of medical pictures and transcript data is done.

## 2. DEEP LEARNING IN MEDICAL DATA ANALYSIS
### A. DIFFICULTIES IN MEDICAL DATA ANALYSIS WITH DEEP LEARNING

Deep learning has always shown the remarkable achievements in speech recognition, scene recognition it enabled the researchers in medical data analysis to embark the milestone, however there are challenges in the field of analysing the medical data. Some are

- **Aspects of medical data.**

For picture information in clinical information, profound learning is generally to prepare some two-layered plane pictures. So, it is important to sample the three-layered clinical pictures and converts them into two-layered picture data, which will not just decrease the image. The resolution also loses a lot of image information. To prevent the losing of image information there are some ways and the sampling method of Roth et al is one among them. This method implies which is mostly in the subsequent while dimensioning, the tested pictures are not generally prepared independently, however, three spatially symmetrical pictures are consolidated into one RGB picture and the first picture data is held as much as could be expected.

- **Number of clinical data**

Deep learning excels in many sectors due of the abundance of learning data available. These learning data's features have more expressive abilities than features collected manually, allowing for better effect to be achieved. As a result, sufficient learning data sets must be available for the deep learning model.

However, in the field of medical data analysis, there is a severe lack of training data, therefore researchers must figure out how to quickly acquire a significant number of tag data and how to utilise the sparse data to train the deep learning model and provide acceptable results.

- **Unsupervised training and supervised training**

The majority of deep learning algorithm-based tasks, including scene identification, picture segmentation, object recognition, etc., all make use of supervised training techniques. However, some research has also been done on unsupervised training techniques, such as image coding and picture representation. The unsupervised training approach is better appropriate for the analysis of medical pictures since many data cannot be used with the deep learning model due to the inherent limitations of medical data.

A Boltzmann machine RBM is part of the unsupervised deep learning model that is currently being used. From the unlabelled training dataset, this method model may directly develop effective feature representations. The drawback is that these characteristics may not always have the optimal categorization impact.

### B. DEEP LEARNING METHOD IN MEDICAL DATA ANALYSIS

Big data in the field of medicine is multimodal, complicated, and holds a plethora of knowledge. The challenges posed by medical big data include how to efficiently gather and obtain medical health data

quickly and accurately, how to use high-speed networks to deliver medical health data reliably and efficiently, how to use machine learning and deep learning techniques related to artificial intelligence to extract useful information from medical big data, and how to develop intelligent applications for the majority of medical staff and regular people.

The artificial neural network model serves as the foundation for the deep learning technique. By integrating numerous nonlinear processing layers, the raw data is abstracted layer by layer. Various degrees of abstract characteristics are then extracted from the data and utilised for classification prediction. Comparing it to conventional machine learning techniques reveals the following three features:

- Architecture for deep models. The deep learning model's multi-layer architecture closely resembles the visual processing system of an animal. Deep learning models include more hidden layers and more nonlinear transformations than other shallow models, such the support vector machine, which considerably improves their capacity to fit complicated models.

- Depiction of data features in layers. The deep learning model summarises the higher-level feature representation to explain the complicated data structure by first taking the original form of the data as input and then using the output of the current layer as the input of the next layer, stacking layer by layer.

- Unsupervised instruction. The training impact is significantly enhanced by the deep learning model's addition of an unsupervised learning process, which it acquires through pre-training. Additionally, training with unlabelled data expands the amount of data that is readily available.
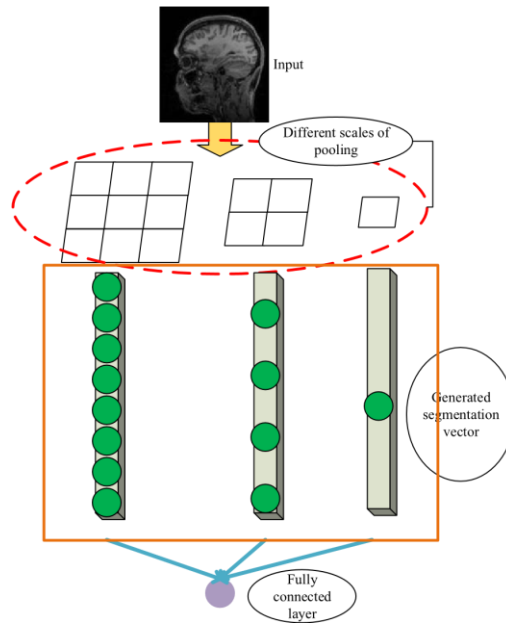
## 3. DEEP LEARNING IN MEDICAL DATA ANALYSIS

*Medical text data analysis model based on 3d convolutional neural network*

The 2D convolutional neural network's limitations, which include its ability to exclusively analyse two-dimensional static input, are addressed by the 3D convolutional neural network. It excels at processing tasks including motion identification with temporal information, natural language processing, and video processing. In order to extract information between the phrase internal information and the phrase of the input data, this article employs a 3D convolution kernel. Additionally, the same network topology may be utilised for processing regardless of the volume of incoming data, considerably extending the model's range of possible applications. The distinction between the 2D convolution and its fundamental structure.

A picture serves as the left input for the 2D convolution structure, several images serve as the left input for the 3D convolution structure, and the feature map produced by convolution serves as the right input. The graphic shows how the 2D convolutional neural network scans the input image using a fixed-size convolution kernel, calculates the value after each convolution, and then synthesises the feature map. In order to extract image features, it is essential for the 3D convolutional neural network to maintain the variation characteristics of many pictures across time. In order to create a feature map, the data from sequential photos taken at the same location are first merged into a three-dimensional matrix. This three-dimensional matrix is then convolved using a three-dimensional convolution kernel, and feature values are produced. A feature map incorporates information from many pictures at once, making it possible to analyse three-dimensional data and extract time series information from the incoming data.

The input information and the convolution kernel are both three-dimensional matrices in the 3D convolutional layer. The size of the input data specified here is $Xp \in Rl*k* v$ for ease of explanation, where l is the number of phrases that make up the patient's input data. Among them, v is the length of each word2vec vector and k is the length of a phrase, or the number of embedded vectors included in a phrase. The vector length v in this study is 50. There is just one vector between each phrase when the data for each patient is divided into sentences.

The issue of various feature map sizes is resolved by the space pyramid pooling method. It also extracts features from various sized receptive fields while working under the assumption that the key aspects of feature maps should be preserved. Currently, it is the finest option for processing input data of various scales. Deep learning, when used with any deep learning model, is crucial for handling erratic data. The architectural design of the spatial pyramid pooling method.

## 4. EXPERIMENTAL RESULTS AND ANALYSIS

In order to verify the validity of the proposed method for medical image data, this paper selects the brain MRI medical image data set as medical image data. This paper selects 91 MCI data samples and 79 normal control data samples. Subjects in the normal control group ranged in age from 63.2 to 88.3, and MCI patients ranged in age from 66.5 to 87.3. There was no significant difference in age and gender between the two groups. All subjects were from the database of ADNI. In this paper, the patient's electronic medical record data is selected as the medical texts.

In order to more intuitively represent the difference in data analysis accuracy of different algorithms, this paper draws the fluctuations of the experiment with the accuracy under different input data and different models, as shown in Figure 1 and Figure 2. Figure 1 shows the accuracy of each algorithm as a function of the number of experiments when the input data is medical image time series data.
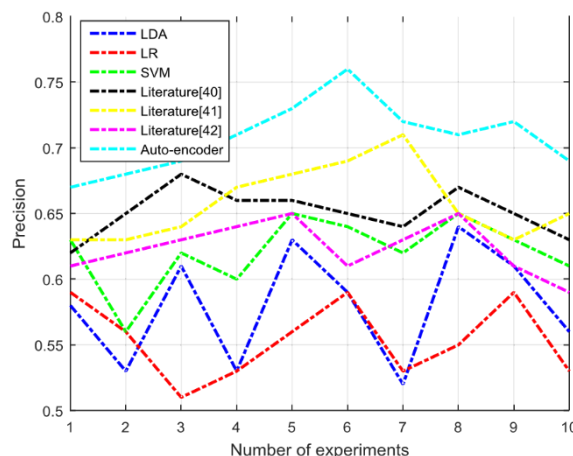


**Figure 1.** The input is the test accuracy of the time series data.

In this experiment, 10-fold cross-validation of all training data was used to identify the ideal AutoEncder network topology. This research first establishes an AutoEncder network with various hidden layer counts, after which repeated tests are carried out for each network layout. In the experiment, it was discovered that once there had been more than ten times as many tests, the average accuracy value would eventually settle and there would be no significant oscillations. Because of this, this research compares the input when the input includes data on correlation coefficients and repeats 10 classification trials for various network designs.

According to Figure 1, if the test's average accuracy is higher, the categorization impact will be more favourable. Figure 1 shows that in 6 out of the 10 tests, the AutoEncder network with three hidden layers had the greatest classification accuracy. Additionally, the two hidden layers' average classification accuracy on the network is 86.47%, the one hidden layer's average classification accuracy on the network is 83.53%, and the three hidden layers' average classification accuracy on the network is 84.70%. In the cross-validation test, the network design with three hidden layers had the greatest classification accuracy. In order to create the analysis model for medical imaging data, this study chooses 3 hidden layers from the AutoEncder network structure.

It's crucial to train the AutoEncder model with the appropriate amount of hidden layer nodes. The data of buried layer nodes is often established empirically in the majority of texts. The network won't be able to fit complicated data if this quantity is too little, and if it's too big, it will lengthen training time and result in over-fitting issues. In this article, we employ a three-layer AutoEncder network with 200 nodes per layer. This research analyses the convergence of the loss function of time series data and correlation coefficient data of medical pictures in order to explain the benefits of employing correlation coefficient data. Figure 2 illustrates how, throughout the fine-tuning phase, the training error loss varies with the number of iterations. As demonstrated in Figure 2, the training error is rapidly decreased at first when there are two different types of input data; however, after about 50 iterations, the training error's convergence speed progressively decreases and tends to be horizontal. At the same time, the training error for the input time series data stabilises at around 510-3 after 100 iterations, whereas for the input data for the correlation coefficient, the training error stabilises at about 110-3 after about 50 iterations.
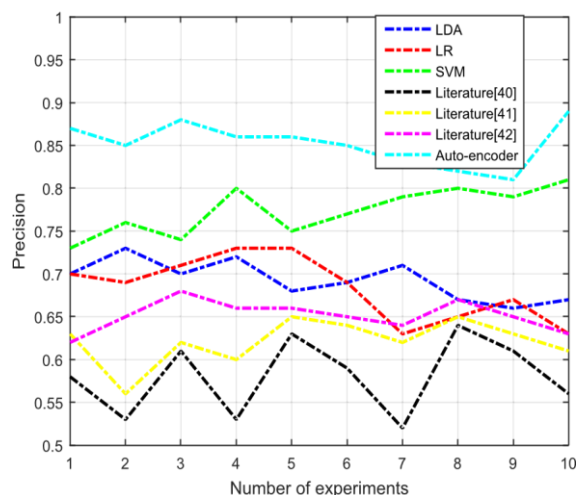


**Figure 2**. The input is the test accuracy of the correlation coefficient data.

## CONCLUSION

The rapid growth of medical big data places more demands on medical big data processing techniques, and the complexity and variety of its forms make analysis more challenging. Deep learning is unquestionably the most successful data processing technique now in use. This study

constructs the related deep learning model independently for the two primary types of medical data—medical picture data and medical text data—and achieves intelligent identification of precise early diagnosis and risk prediction for certain diseases. To begin with, the Auto Encoder deep learning model is created, and it has the ability to pre-train the network in advance and use less time and resources.

The approach is also readily adaptable to other medical picture data analysis and processing, which is crucial for enhancing the precision of illness detection. Second, a deep learning model incorporating spatial pyramid pooling and 3D convolutional neural networks is developed. The spatial pyramid pooling structure can process input data of any length, while the 3D convolution structure in the model can process input data of arbitrary length while extracting the internal features of the data. This allows for effective data analysis for the patient's future in the intelligent identification and control of disease risk.

## References

1. Q. Dou, H. Chen, L. Yu, L. Zhao, J. Qin, D. Wang, V. C. Mok, L. Shi, and P.-A. Heng, ''Automatic detection of cerebral microbleeds from MR images via 3D convolutional neural networks,'' IEEE Trans. Med. Imag., vol. 35, no. 5, pp. 1182–1195, May 2016

2. G. van Tulder and M. de Bruijne, ''Combining generative and discriminative representation learning for lung CT analysis with convolutional restricted Boltzmann machines,'' IEEE Trans. Med. Imag., vol. 35, no. 5, pp. 1262–1272, May 2016

3. M. J. J. P. van Grinsven, B. van Ginneken, C. B. Hoyng, T. Theelen and C. I. Sánchez, ''Fast convolutional neural network training using selective data sampling: Application to hemorrhage detection in color fundus images,'' IEEE Trans. Med. Imag., vol. 35, no. 5, pp. 1273–1284, May 2016.

4. D. Yu and L. Deng, ''Deep learning and its applications to signal and information processing [exploratory DSP],'' IEEE Signal Process. Mag., vol. 28, no. 1, pp. 145–154, Jan. 2011.

5. J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei, ''ImageNet A large-scale hierarchical image database,'' in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., Jun. 2009, pp. 248–255.

6. A. Krizhevsky, I. Sutskever, and G. E. Hinton, ''ImageNet classificatio with deep convolutional neural networks,'' in Proc. Adv. Neural Inf. Process. Syst., 2012, pp. 1097–1105

7. Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick S. Guadarrama, and T. Darrell, ''Caffe: Convolutional architecture for fast feature embedding,'' in Proc. 22nd ACM Int. Conf. Multimedia, Nov. 2014, pp. 675–678.

8. C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, ''Going deeper with convolutions,'' in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., Jun. 2015, pp. 1–9.

9. K. Shameer, M. A. Badgeley, R. Miotto, B. S. Glicksberg, J. W. Morgan, and J. T. Dudley, ''Translational bioinformatics in the era of real-time biomedical, health care and wellness data streams,'' Briefings Bioinf., vol. 18, no. 1, pp. 105–124, Feb. 2016.

10. J. Ma, Y. Yang, D. Huo, Z. Wang, X. Zhai, J. Chen, H. Sun, W. An, J. Jie, and P. Yang, ''LincRNA-RoR/miR-145 promote invasion and metastasis in triple-negative breast cancer via targeting MUC1,'' Biochem. Biophys. Res. Commun., vol. 500, no. 3, pp. 614–620, Jun. 2018.

11. Z. Li, J. Li, A. Ren, R. Cai, C. Ding, X. Qian, J. Draper, B. Yuan, J. ng, Q. Qiu, and Y. Wang, ''HEIF: Highly efficient stochastic computing-based inference framework for deep neural networks,'' IEEE Trans. Comput.-Aided Design Integr. Circuits Syst., vol. 38, no. 8, pp. 1543–1556, Aug. 2019.

12. M. Eberl, D. Mangelberger, J. B. Swanson, M. E. Verhaegen, P. W. Harms, M. L. Frohm, A. A. Dlugosz, and S. Y. Wong, ''Tumor architecture and notch signaling modulate drug response in basal cell carcinoma,'' Cancer Cell, vol. 33, no. 2, pp. 229–243, Feb. 2018.

13. F. Alfaro-Almagro et al., ''Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank,'' Neuroimage vol. 166, pp. 400–424, Feb. 2018

**14.** I. Tatsuya, ''Information processing unit,'' J. Mater. Process. Technol.,vol. 3, no. 8, pp. 165–196, 2018.

**15.** D. Sumathi and P. Poongodi, ''Secure medical information processing in cloud: Trust with swarm based scheduling,'' J. Med. Imag. Health Informat., vol. 6, no. 7, pp. 1636–1640, Nov. 2016.

**16.** G. S. Birkhead, M. Klompas, and N. R. Shah, ''Uses of electronic health records for public health surveillance to advance public health,'' Annu. Rev. Public Health, vol. 36, no. 1, pp. 345–359, Mar. 2015.

**17.** L. C. Gurrin, J. J. Kurinczuk, and P. R. Burton, ''Bayesian statistics in medical research: An intuitive alternative to conventional data analysis,''J. Eval. Clin. Pract., vol. 6, no. 2, pp. 193–204, May 2010.

**18.** A. M. Hassan, Y. F. Hassan, and M. H. Kholief, ''A deep classification system for medical data analysis,'' J. Med. Imag. Health Informat., vol. 8, no. 2, pp. 250–256, Feb. 2018.

**19.** K. Lu, ''Number of imputations needed to stabilize estimated treatment difference in longitudinal data analysis,'' Stat. Methods Med. Res., vol. 26, no. 2, pp. 674–690, Oct. 201

**20.** V. Taslimitehrani, G. Dong, N. L. Pereira, M. Panahiazar, and J. Pathak,'Developing EHR-driven heart failure risk prediction models using CPXR(Log) with the probabilistic loss function,'' J. Biomed. Informat., vol. 60, pp. 260–269, Apr. 2016.

**21.** Z. P. Liu, L. Y. Wu, Y. Wang, X. S. Zhang, and L. Chen, ''Prediction of protein-RNA binding sites by a random forest method with combined features,'' Bioinformatics, vol. 26, no. 13, pp. 1616–1622, 2010.

**22.** J. H. Chen and J. Z. Lin, ''Developing an SVM based risk hedging prediction model for construction material suppliers,'' Autom. Construct., vol. 19, no. 6, pp. 702–708, Oct. 2010.