

# Assisting clinicians in predicting affected livers from patient data using DLAPDL

<sup>1</sup>Dr. R. MALATHI

<sup>1</sup>Assistant Professor, P.G. and Research Department of Computer Science, H.H. The Rajah's College (A), Pudukkottai, shamalu2010@gmail.com

## Abstract

Liver, biggest organ in the human body, controls most metabolic activities in humans including converting nutrients into food, bile synthesis, protein creation, glucose storage/, processing cleaning the blood, immunological component production, bilirubin clearance etc. Thus, it is an important and crucial body organ, and its maintenance is primary to human health. Liver is overlooked by humans due to their unhealthy lifestyle routines, thus resulting in acute to severe liver problems like LCs (Liver Cancers). Healthcare systems have been using automations in health related decisions where the systems extract relevant information from massive medical datasets using MLTs (Machine Learning Techniques) which assist clinicians in taking accurate and quick choices in terms of illness predictions or diagnosis. This paper proposes an automated diagnostic framework called DLAPDL- (Deep Learning Approach for Prediction of Diseased Liver) based on CNNs (Convolution Neural Networks) to estimate DLs (Diseased Livers) based on medical reports data of patients. In this respect, this paper presents an in-depth examination of predictions of DLs. The proposed approach's results show higher levels of accuracy in experimentations and evaluations on medical datasets.

**Keywords:** Machine Learning, Diseased Livers, Deep Learning, Exploratory Data Analysis, CNNs.

## INTRODUCTION

DLs in humans have been growing as a result of physical inactiveness where urban and metropolitan areas report increasing evidences. DLs are caused by many factors including viruses (hepatitis A/B/C), immune systems issues like autoimmune hepatitis, primary biliary/sclerosing cholangitis). Drugs, poisons, or high alcohol consumptions can also lead to DLs like fatty DLs and non-alcoholic DLs, non-alcoholic steatohepatitis, cirrhosis or inherited diseases (hemochromatosis, hyperoxaluria, alpha-1 antitrypsin deficiency, and Wilson disease) and Cancers including Liver, bile ducts and cell adenoma). Millions of people die due to DLs every year where viral hepatitis kills 1.34 million people perennially. India is at a higher risk of DLs and may

become DLs World Capital by 2025. This is mainly due to the deskbound lifestyles, increasing use alcohol/smoking culminating in DLs. There are over a hundred different kinds of liver infections and it is imperative to be concerned about combating these illnesses. Figure 1 depicts DLs.

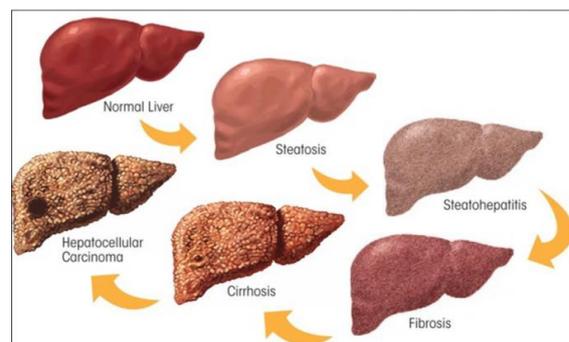


Fig. 1 – Diseased Livers

One of the primary tasks of the liver is to cleanse the blood of poisons present in the food. Inflammations are the initial symptoms of DLs. Inflammations in the liver indicate that human immune system is reacting to external substances like poisons. They are caused by a variety of factors like tenderness and swelling due to higher level of fat in the liver or too many toxins or viral infections like Nonalcoholic fatty liver and liver hepatitis and viral hepatitis result in inflamed livers. Scarred liver tissues do not function the same way as healthy liver tissues, resulting in fibrosis. Scarred connective tissues obstruct blood flow into the liver resulting in dysfunctions. Toxins can accumulate in the brain, causing problems with attention, memory, sleeping, and other mental processes. Once diagnosed as DLs, treatments focus on preventing worsening of the liver by protecting healthy liver tissues. Cirrhosis is a condition where liver's exhausting connective tissues are scarred and as it progresses without treatment, liver fails to function. Cirrhosis cause liver disorders or DLs and are one of the primary causes for LCs, the other being Hepatitis B. Carcinomas can develop in the liver any time while the disease is progressing. Blood tests, blood counts, abdominal and pelvic CT scans/ultrasounds, isoenzyme testing, lactate dehydrogenase tests, liver biopsies, liver function tests and Needle Biopsies are all used to assess the condition of liver. In spite of these exhaustive procedures found for diagnosing DLs, conventional methods of treating DLs face limitations. Healthcare images or diagnostic procedures generate large amounts of medical data is produced compared to expert analysis of significant data. Clinicians without much expertise find it complex to analyze healthcare images. Multiple issues exist in the case of DL diagnosis. Though uncovering hidden patterns/relationships from images is pertinent to data processing, images may contain noises or be of poor quality due to issues in imaging modalities, medical image processing gets affected in diagnosis. Liver Biopsies are considered dangerous and are clinician dependent and may result in inaccuracies of predictions. One biomarker is not sufficient for predicting DLs as they require many

biomarkers, making manual detections sluggish and confusing. Moreover, scarcity of liver for transplants adds to existing complexities. Thus, prevention is better than cure in the case of liver malfunctions where its early diagnosis can save lives resulting in reduction of mortality rates around the globe. To address the limitations of traditional illness diagnosing methods, automatic systems can be utilized. Clinical choices are made with the primary goal of providing an accurate and fast diagnosis. As a result, an automated disease detection system can assist a physician or even a patient in making precise and accurate DLs predictions. Clinical decision-making software [1] assist clinicians and patients in making prompt diagnoses and receiving required therapy to reduce illness consequences. The use of MLTs in these systems has been proposed in many studies. As advances in computer science are accelerating DLNs are gaining popularity for predictions in a variety of disciplines, this study proposes a DLN approach for predicting DLs from patient information. The scheme DLAPDL categorizes a patient as DL or not DL using CNNs. This introductory section is followed by a study of related literature while its subsequent section details on the proposed approach. The fourth section is results and discussions followed by conclusion which the fifth and final section. Review of DLs Studies: For a long time, AIs (Artificial Intelligence) have been utilized to identify LDs and analyze patient's data for inferences on early illness diagnostics, and therefore aiding clinicians and patients in determining treatments. KDDs (Knowledge discovery in Databases) is also an evolving branch of computer science and engineering that is commonly utilised for extracting crucial and significant information from large medical datasets. These inferences are needed by medical practitioners where they have to provided within short durations. Prior studies could point out that bilirubin was the most significant element in liver assessments. The study in [2] predicted PSCs (Primary Sclerosing Cholangitis) in [2] using MLTs with enhanced outcomes when compared to C-Statistics and enabled patients to self-access online. Non-invasive MLTs diagnosed anthropometric and serum biomarkers in [3].

The study's approach used NNs (Neural Networks) was a cheap and easy way as it modified SAFs (Steatosis, Activity, and Fibrosis) for evaluating presence of DLs. The algorithm used biomarkers produced at different stages of DLs in Southeast Asian country in patient information for training. The changes were detected in serum for predicting DLs. The study in [4] introduced DNNs (Deep NNs) for screening of DLs with a four-layer DNN. The study used fifteen characteristics to screen 76,914 patients. Using the same dataset, the model's prediction performance was compared with LRs (Logistic Regressions), RFs (Random Forests) using 5-fold cross-validations. The study's proposed DNN model outperformed other models in comparisons with modest over-fitting in the screening process. Disease-free survival of patients after liver resections were assessed in [5] using RFs which used clinical and other parameters. The inputs were pre-processed and trained by RFs, for identifying low/ high risk patients. The study also used RFEs (Recursive Feature Eliminations) for selecting relevant predictors and thus facilitated their model's accuracy. The study helped identify line of treatments to DL patients based on the risk factor. The study in [6] proposed a non-invasive way to assess HCCs (Hepatocellular Carcinoma). by evaluating 4423 chronic hepatitis C patient records. The study identified significant parameters for HCC diagnosis and validated their scheme using a dataset with records of 293 HCV, 53 HCC patients. The model's deployment using multitude of MLTs showed that their alternating DTs (Decision Trees) showed best results with an ROC (Receiver Operating Characteristic) curve value of 95.6% and an accuracy of 99%. The study projected the importance of albumin and TB parameters in HCC predictions. Based on breath-based metabolites like VOCs (Volatile Organic Compounds), the study in [7] suggested non invasive approaches for determining patient's health in terms of pulmonary hypertensions, cirrhosis and HCCs. The study's RFs learnt from 22 VOCs of 296 patients while their findings reported higher levels of two VOCs namely Acetaldehydes and Acetones in cirrhosis and HCC affected patients. The model

was an aid for accurate detection of illnesses and tumors in DLs. MRIs (Magnetic Resource Imaging) were used to detect HCC elements (microvascular invasions - mVI). The study used MLTs before surgeries or in a line of therapies. By manually segmenting the HCC region, radiomic characteristics were retrieved while MLTs were trained on characteristics obtained from MRI sequences of tumour or peritumoral regions or a combination of both regions. The study showed accuracies above 86% in all three types of regions and this could forecast the presence of mVIs in the lungs. The study [8] suggested skewing data for MLTs and predicted postoperative survivals after transplantations of the liver. Their suggested scheme, MELD (Model for End-stage Liver Disease) scored data as a pre-processing function while 10-fold cross-validation RFs selected features and categorized data. The suggested model when compared with LRs and DTs showed that RFs provided higher accuracy (77.1%) in predicting INRs (International Normalized Ratios), lymphocytes, platelets, WBCs (White Blood Cells), Mg (magnesium) and Na (sodium). The use of a MLT based radiomic model in the prediction of metachronous metastases of CRC (colorectal cancer) patients was examined in [9]. The dataset had Computed Tomography images of 91 patients. The study split it into two groups based on the presence of metachronous liver metastases. The scheme processed and segmented images yielding 1767 radiomic characteristics for each patient and eliminating inter-correlated elements using the Kruskal-Wallis test. Their use of RFs, wrapper based feature selections and Bayesian optimizations were built utilising clinical data, radiomic data, and a mix of both clinical and radiomic data. The AUC for three models was 86% (with radiomic data), 71% (with clinical data), and 86% (without any data). Their radiomic data analysis yielded important characteristics like biomarkers, for identifying individuals with high risk of developing colorectal liver metastases. DeepProg, proposed in [10] was a unique and general computational framework for predicting survival by analyzing different types of data sets using a combination of DLNs and MLTs. DeepProg outperformed other

multi-omics integration technique in terms of prediction accuracy. Though many studies have reported DLs in their schemes while using patient information as inputs, gaps do exist in terms of accuracy. This research work attempts to improve the accuracy of DL predictions using DLNs.

**DLAPDL:** Humans require proper healthcare in their current lives. There is a need for medical services that are easily available to everyone. The proposed work is implemented on IPLDs (Indian Liver Patient Datasets) of Irvine database from University of California. The dataset encompasses several attributes which are used by the proposed system for DL predictions. The main focus of the proposed work is to predict DLs based on DMTs and software engineering approaches where data is prepared, attributes are analyzed and finally DLs are classified using patient's information from the dataset. Records of patients help in the prognosis of mild to severe DLs as they include different test results. The dataset had ten variables related to health where the output of the proposed system was in a binary format DLs or non-DLs. Figure 2 depicts the proposed DLAPDL framework.

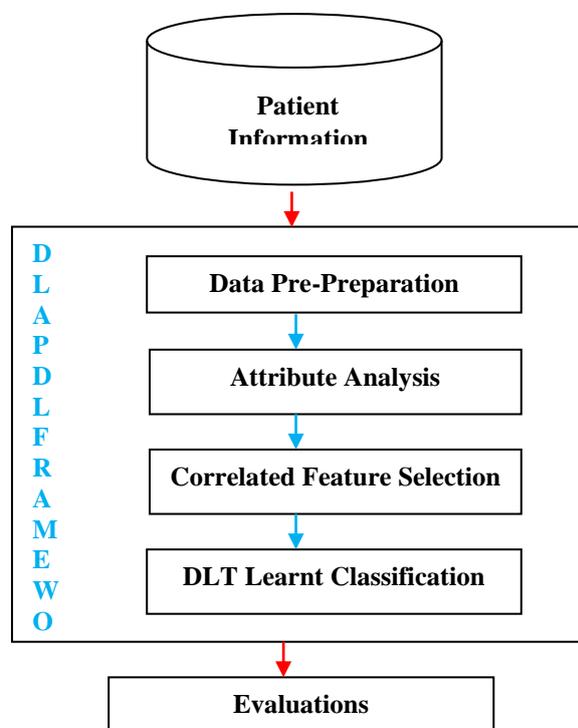


Fig. 2 - DLAPDL Framework

**DLAPDL Data Preparation:** Data preparation is an important part of processing as the data needs to be converted into a form that can be used for processing or inferring required information. Since, impurities or differentiable values can exist in data leading to less accurate results. Though there are multiple preparation techniques, they generally aim to solve deficiencies in input data. In the case of the used dataset, only null values needed to be normalized and hence data was prepared based on this factor for further processing. This is subsequently followed by attribute analysis.

**DLAPDL Attribute Analysis:** Data scientists usually employ exploratory analysis for ensuring accuracy of results mainly for acceptance by business establishments. It benefits stakeholders as it paves the way for proper answers to their questions. They prelude to groupings or reductions of features and assist these techniques by: Graphical descriptions of high dimensional data; summarized statistics; univariate displays of raw data values/dataset's unprocessed fields; bi-variate visualisations and summarization of statistics for links between each variables in the dataset; multi-variate visualizations which use statistics and data to map and explain relationships between multiple fields in data and predictive models including MLTs to predict outcomes. The attributes of the dataset are detailed below:

- **Total\_Bilirubin:** Bilirubin, an orange-yellow pigment breaks down a portion of red blood cells. The quantity of bilirubin in blood is measured by a bilirubin test and finds out health issues causes like DLs. Bilirubin in the blood is passed as bile by the liver. When bilirubin levels are greater than normal, it implies liver isn't efficiently breaking down wastes.
- **Direct\_Bilirubin** –Direct/conjugated bilirubin is linked to glucuronic acid derived by the liver and refers to the total amount of bilirubin in the blood.
- **Alkaline\_Phosphatase** –Alkaline phosphatases are blood enzymes that aid in breaking down proteins. It has a variety of

functions, but is crucial for functioning of the liver and bone formation. This is assed using ALP tests.

- **Alamine\_Aminotransferase** –Alanine aminotransferase is predominantly found in the liver and kidney. Serum has a low amount of Alanine aminotransferase and helps in monitor DLs. An increased value implies the liver is damaged.
- **Aspartate\_Aminotransferase:** The enzyme aspartate aminotransferase is mostly located in the liver, although it can also be found in the muscles. It is released into circulation when liver is damaged and measured by an AST blood test. Its value is used by clinicians to determine DLs or illness.
- **Total\_protein** – The human body has two kinds of protein: albumin and globulin. total protein test value determines their normal health. Its value is assessed in cases of sudden weight loss or tiredness or kidney or DLs symptoms.
- **Albumin** - Liver produces albumin which maintains fluid in circulation of Hormones, vitamins, and enzymes while preventing its leakage into other tissues. Low albumin levels indicate renal diseases or DLs.

**Albumin\_and\_Globulin\_Ratio:** The Albumin to Globulin ratio is the ratio of albumin present in serum in relation to the amount of globulin and interpreted as total protein concentrations. Its default ratio is 1:1. Figure 3 depicts the Albumin Ration of Patients.

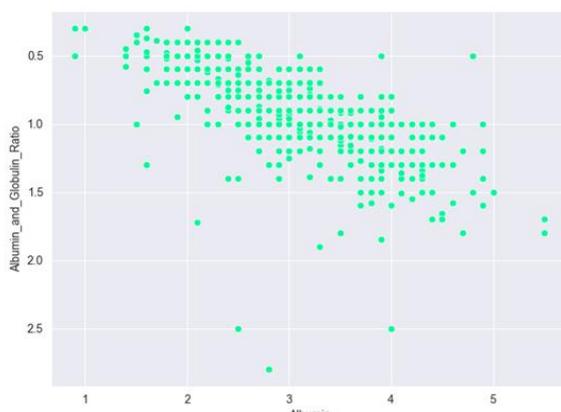


Fig. 3 – Graphed Albumin Ration of Patients

SDs (Standard deviations), categorical variables, and confidence intervals are all used in attribute analysis and thus insights of data are derived. The analysis characteristics can be used for further advanced data analysis or modeling. Figure 4 depicts Attribute Analysis Output.

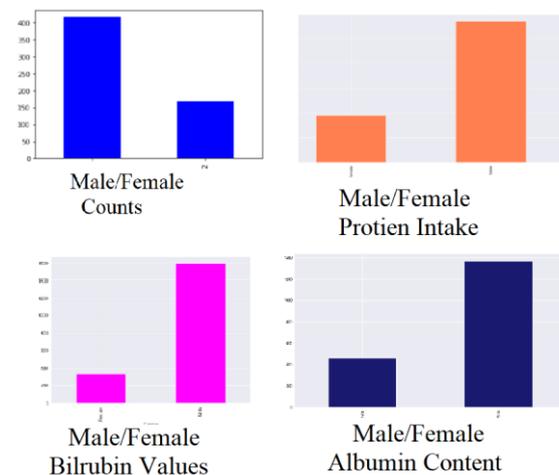


Fig. 4- DLAPDL Attribute Analysis

**DLAPDL Feature Identification:** Feature identification is the selection of relevant features that imply strong relations of features to outputs, assisting in faster learning by models or minimizing dimensionalities/complexities, with proper interpretations resulting in enhanced accuracy of predictions. Features can selected using Univariate identifications or based on their importance or correlations. Univariate identifications are statistical tests (P-Value or Chi Square Test) used to select features with strongest relationship with the output variable. Feature Importance is assessing the importance of each feature in the dataset. Feature importances are scores given for each feature and higher score features are important. Correlations show feature relationships amongst features or target variables which can be positive (proportional), negative (Inversely Proportional). Figure 5 depicts DLAPDL attribute Correlations.

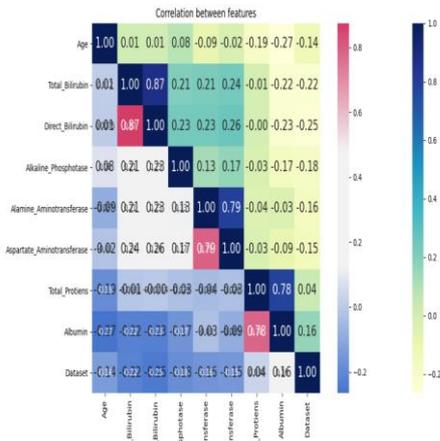


Fig. 5 – DLAPDL Attribute Correlations

DLAPDL Classification: Classification is an important use of DMTs and used for differentiating classes found in data and based on their learning from training data (objects whose classes are unknown). Automatic recognition of complex patterns and making intelligent decisions based is a predominant area of research. Classification can also be called as Supervised learning. This work uses DLTs and specifically CNNs (Convolution Neural Networks) for classification of DLs. The proposed DLAPDL uses CNNs to classify test data into two classes namely DFLs and non-DLs. CNNs have input and output layers with multiple hidden layers which encompass convolution, pooling and fully connected layers. Convolution (simulating the response of a neuron) operations on the inputs occur in the convolution layer which is then passed to the next layer. CNN architecture used in this study is depicted as Figure 6.

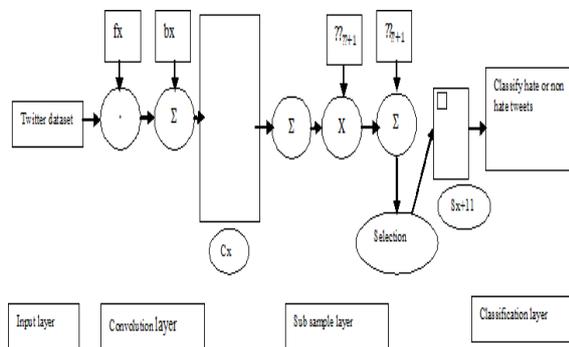


Fig 6 - Architecture of Modified CNN

CNNs may incorporate local or global pooling layers which join neuron’s cluster outputs into the neuron of the subsequent layer. Mean pooling is implementing neuron cluster average values in the previous layer while Fully connected layers map each neurons of layers to neurons in other layers. CNNs are principled like traditional MLP (Multi-Layer Perceptron) neural networks [11]. CNNs have input/convolution/ sub-sampling/ classification layers and can analyze high-dimensional data efficiently. Parameter sharing of convolved layers reduce parameter counts. Input layer gets its input from samples which are then transformed and submitted to the next layer. Initial parameters like local fields scale and filters are also defined in this layer. Cx (Convolution layer) convolutes inputs for producing several layers called feature maps which include previous convolution layer computations. This layer extracts key features and reduces the computational complexity of the network. An activation function is executed in the convolution layer. The function maps outputs to a set of inputs creating a non-linear network structure. Weights are added to feature values for a new pattern output defined as Equations (1) and (2)

$$y(n) = f(\sum_{i=1}^{i=N} w_i(n)x_i(n)) \tag{1}$$

$$\text{Where } f(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases} \tag{2}$$

Where n stands for the iteration index

The Connection weights are then updated based on (3)

$$w_i(n + 1) = w_i(n) + \eta(d(n) - y(n))x_i(n), \quad i = 1, 2, \dots, N \tag{3}$$

Where η is the gain factor

And the application of SD is defined in Equation (4)

$$\sigma = \sqrt{\frac{1}{n} \sum f_i(x_i - \bar{x})^2} \tag{4}$$

The weighted features are fed to CNN for better classification accuracy. A Polynomial distribution function ensures that the same set of data is analyzed. Every feature map generated in a convolution layer is sub-sampled in this layer.

**Results and Discussion:**

Stage wise experimental results of Python 3.9 on windows 10 using 4 GBis displayed in this section. The code was targeted on ILPD

datasetwith data on 10 DLs patients. The dataset encompasses 583 records with 416 DLs apart from 167 patients who do not have a liver disease. Figure 7 depicts a snapshot of the dataset.

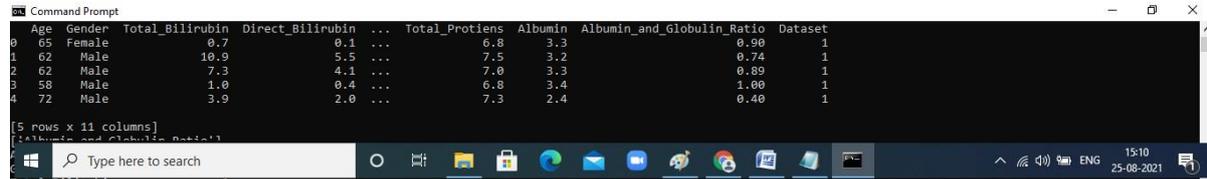


Fig. 7 Snapshot of the IPLD dataset

DLAPDL Data Preparation Results: The dataset was in CSV format had multiple attributes with demographic and health related information. The proposed framework prepared input data by checking for null values. The

parameter 'Albumin\_and\_Globulin\_Ratio' had null values and hence, all null columns were replaced with a mean value of 0.5. Figure 8 depicts the output of DLAPDL Data Preparation.

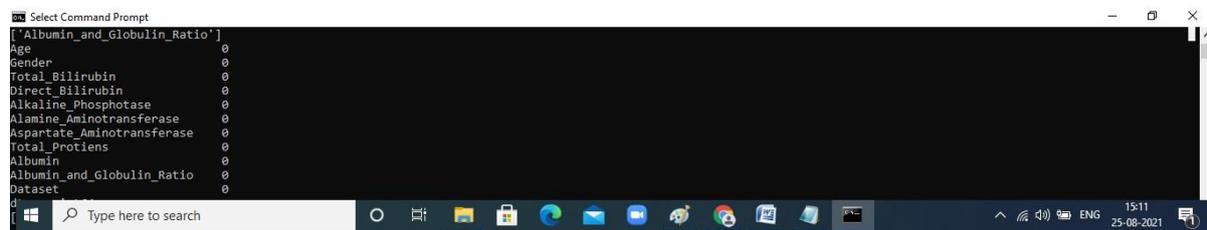


Fig. 8 - DLAPDL Data Preparation Output

DLAPDL Attribute Analysis Results: The main purpose of attribute analysis is assessment of data before making any assumptions. It can help discovering obvious errors while promoting better understanding of data trends, outlier identification, unusual event detections and discovery of intriguing relationships between variables. Non-numerical columns like gender were made numeric where males were assigned 1 and females Zero. The number of male and females were counted. The dataset had more males than females. Protein intake was found to be higher in Males. The level of Albumin was also high in Males. Further, males also had higher Bilirubin contents when compared to females. It should be noted that higher Bilirubin content implies higher DLs. Figure 9 depicts DLAPDL's pair wise Attribute Analysis.

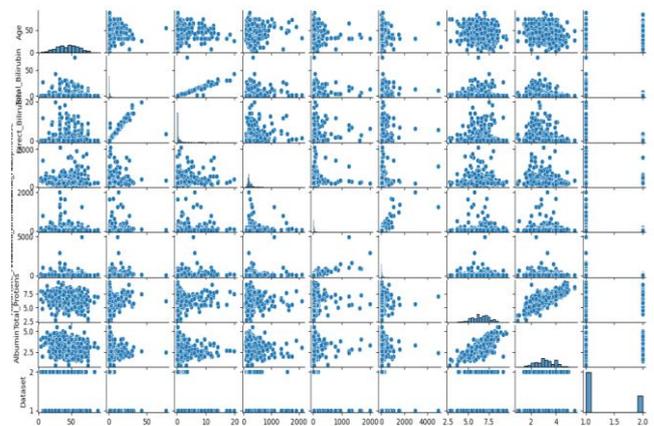


Fig. 9 - DLAPDL's Attribute Analysis Output

DLAPDL Feature Identification: Statistical approaches ad MLTs extracted significant biomarkers from patient records in the ILPD dataset. Gain ratios compared the usefulness and ranks of features, thus demonstrating the model's dependability. This work used positive correlations between features and the attribute's mean values to understand corrections better and thus unwanted features wee eliminated. Figure 10 depicts the output of DLAPDL feature mean values and correlations.

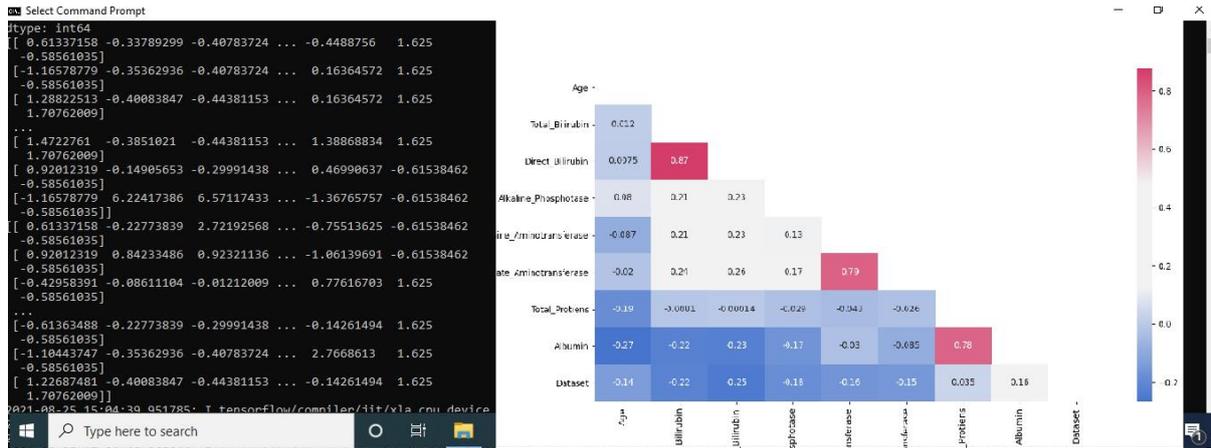


Fig. 10 - DLAPDL Feature Identification Output

DLAPDL Classification: The proposed DLAPDL improves on classification by its use of DLTs with enhanced polynomial distributions. The DLT model was used to predict the class label of objects for which class

label of DLs and non-DLs is unknown. DLTs proposed in the study learn (or improve their performance) based on IPLD data. Figure 11 depicts DLAPDL Training for Precitions.

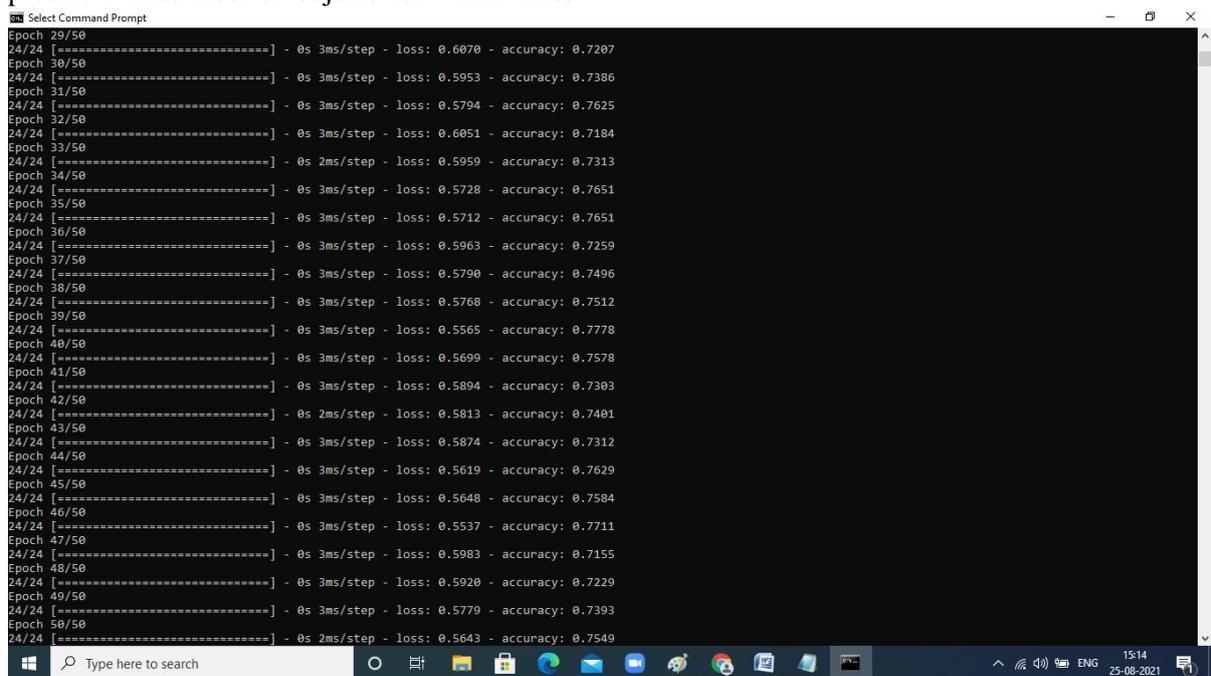


Fig. 11 DLAPDL Training Screen shot

Data Set Splits and Cross Validations are approaches used in studies where subsets of datasets are separated for a model's testing in cross validations. It helps prevent overfitting/underfitting issues. Hold-out (early breaks) and K-folds are two cross validation approaches used often[2]. Data splits also help eliminate bias where the ratio of splits varies based on the application or model. Examples of sampling include random, trial-and-error,

systematic, convenience, CADEX, DUPLEX sampling [2] where simple random sampling are most frequently used by MLTs in their sampling. MLTs learn and prediction based on training data and validate it on test data. Thus, models learn in training not validations where they improve hyperparameters. The test data is used to assess the performance of the constructed model once the hyperparameters have been determined. The three main metrics



Table 1 – *Comparative Performances of Algorithmic accuracy*

S.NO	DLT	Accuracy
1	Naive Bayes	55.65%
2	SVM	71.38%
3	KNN	71.52%
4	DLAPDL	74.8%

## CONCLUSION AND FUTURE WORKS

With the passage of time, diseases of the liver and heart are becoming increasingly widespread. These are only likely to get worse in the future, thanks to ongoing technology improvements. Despite the fact that people are becoming more health-conscious and enrolling in yoga and dancing classes, the sedentary lifestyle and amenities that are constantly being introduced have continued to be an issue for their resulting ailments. Hence, the application of the suggested DLAPDL in diagnostics can be very beneficial to society and the world in general as the suggested approach achieved 74.8 % in terms of accuracy, a difficult task while working with complex datasets. It can be concluded that the suggested approach can predict risks of DLs with greater accuracy. Today almost everybody above the age of 12 have smartphones which can incorporate this work's suggested approach in addition to websites and thus be beneficial for a large section of society. Moreover, 74.8 % accuracy on the validation set, with further tuning of this model on more validation data, it might be useful as a tool to suggest further definitive testing of LDs for patients. It can also be used in the future for early detection of LDs non-invasively like mammograms, which also is indicative of breast cancer, but not diagnostic. Future would also be in applying this model to diverse data set with more even gender distributions.

## Reference

- [1] Książek W, Abdar M, Acharya U R and Pławiak P 2019 A novel machine learning approach for early detection of hepatocellular carcinoma patients *Cogn. Syst. Res.* 54 116–27
- [2] Shung D L and Assis D N 2020 Machine Learning in a Complex Disease: PREsTo Improves the Prognostication of Primary Sclerosing Cholangitis *Hepatology* 71 8–10
- [3] Aravind A, Bahirvani A G, Quiambao R and Gonzalo T 2020 Machine Learning Technology for Evaluation of Liver Fibrosis, Inflammation Activity and Steatosis (LIVERFASt<sup>TM</sup>) *J. Intell. Learn. Syst. Appl.* 12 31–49
- [4] Yao Z, Li J, Guan Z, Ye Y and Chen Y 2020 Liver disease screening based on densely connected deep neural networks *Neural Networks* 123 299–304
- [5] Schoenberg M B, Bucher J N, Koch D, Börner N, Hesse S, De Toni E N, Seidensticker M, Angele M K, Klein C, Bazhin A V., Werner J and Guba M O 2020 A novel machine learning algorithm to predict disease free survival after resection of hepatocellular carcinoma *Ann. Transl. Med.* 8 434–434
- [6] Hashem S, ElHefnawi M, Habashy S, El-Adawy M, Esmat G, Elakel W, Abdelaziz A O, Nabeel M M, Abdelmaksoud A H, Elbaz T M and Shousha H I 2020 Machine Learning Prediction Models for Diagnosing Hepatocellular Carcinoma with HCV-related Chronic Liver Disease *Comput. Methods Programs Biomed.* 196
- [7] Nebbia G, Zhang Q, Arefan D, Zhao X and Wu S 2020 Pre-operative Microvascular Invasion Prediction Using Multi-parametric Liver MRI Radiomics *J. Digit. Imaging*
- [8] Liu C L, Soong R S, Lee W C, Jiang G W and Lin Y C 2020 Predicting Short-term Survival after Liver Transplantation using Machine Learning *Sci. Rep.* 10 1–10.
- [9] Taghavi M, Trebeschi S, Simões R, Meek D B, Beckers R C J, Lambregts D M J, Verhoef C, Houwers J B, van der Heide U A, Beets-Tan R G H and Maas M 2020 Machine learning-based analysis of CT radiomics model for prediction of

- colorectal metachronous liver metastases  
Abdom. Radiol.
- [10] Poirion O, Chaudhary K, Huang S and Garmire L X 2020 DeepProg: an ensemble of deep-learning and machine-learning models for prognosis prediction using multi-omics
- [11] Sugauma, Masanori, Shinichi Shirakawa, and Tomoharu Nagao. "A genetic programming approach to designing convolutional neural network architectures." Proceedings of the genetic and evolutionary computation conference. 2017.