



Republic of Iraq  
Ministry of Higher  
Education and Scientific  
Research  
University of Diyala



# *Diagnosis Of Liver Disease Using Machine Learning Algorithms*

A Thesis

Submitted to the Department of Computer Science\ College  
of Science\ University of Diyala in a Partial Fulfilment of the  
Requirements for the Degree of Master in Computer Science

*By*

*Sondos Jameel Mukhyber*

*Supervised By*

*Prof. Dr. Dhahir Abdulhade Abdulah  
Prof. Dr. Amer Dawood Majeed*

2021 A.D.

1442 A.H.

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ  
يَرْفَعُ اللَّهُ الَّذِينَ آمَنُوا مِنْكُمْ  
وَالَّذِينَ أُوتُوا الْعِلْمَ دَرَجَاتٍ

صدق الله العظيم

( سورة المجادلة الاية: 11 )

# Acknowledgment

*Firstly, all my prayers go to Allah, the Almighty, who helped me to accomplish this work until its end.*

*I wish to express my thanks to my college (college of science), my supervisors, Prof. Dr. Dhahir Abdulhade Abdulah and Prof. Dr. Amer Dawood Majeed for supervising this research, and for the generosity, patience, and continuous guidance throughout the work. It has been my good fortune to have advice and guidance from them. My thanks to the academic and administrative staff at the Department of the computer sciences.*

*I would like to thank my family who has supported me throughout and encouraging me to further my studies and help to complete this project, and I wish to express my thanks to my friends.*

*Sondos jameel*

# *Dedication*

*I would like to dedicate this work to:*

*Whom taught me that the champions will  
never be defeated, but they convert it to  
victory.*

*Our Prophet Mohammed*

*Peace be Upon Him (PBH)*

*My Father, My Mother, My Sisters, and  
My Brothers.*

*Sondos jameel*

## *Abstract*

In the present time, liver disease becomes one of the significant topics for researchers in the field of medical technology. Patients with liver disease can only be saved if and only if the disease is early discovered, otherwise the chance of survival chances is slim. The most significant issue, according to this, is early and accurate detection, which is crucial for curing this disease. A Classification is one of the main issues in the knowledge discovery fields and sciences decision. There are many types of algorithms used for constructing classifiers and liver disorder detection.

In this thesis, efficient liver disease classification model have been built to increase the accuracy and decrease the error rate in the disorder detection process. The proposed model consists of four main stages: pre-processing, split of data, z- score normalization and classification of the liver disease stage. Datasets may contain some missing values. In this proposed model the way to handle these values is done by replacing each missing value with the mean value. In this thesis, five algorithm (Decision Tree, Random Forest, Artificial Neural Networks, Support Vector Machine, and K-Nearest Neighbour) are used.

The proposed model has been tested by using two liver datasets (Indian Liver patients and Iraqi liver patients). In comparison to other existing approaches, the results show that the model has a high accuracy rate, the accuracy of the Decision Tree, Random Forest, Artificial Neural Networks, Support Vector Machine, and K-Nearest Neighbour using Iraqi liver patient dataset is 99.06%, 99.06%, 95.32%, 89.71, and 87.85% and when using Indian Liver patient Dataset the accuracy rate is 80.34%, 77.77%, 84.61%, 80.34%, and 77.77%.

## List of Contents

<i>Contents</i>	<i>Page No</i>
<b><i>Chapter One: General Introduction</i></b>	
1.1 Introduction	1
1.2 Related Works	2
1.3 Problem Statement	4
1.4 Aim of Thesis	5
1.5 Outline of Thesis	5
<b><i>Chapter Two: Theoretical Background</i></b>	
2.1 Introduction	6
2.2 An Overview of Liver Diseases	6
2.3 Liver Functions Test	7
2.4 Knowledge Discovery Process in Database	8
2.5 Data Mining Applications in Healthcare	10
2.6 Data Normalization	11
2.6.1 Z-score Normalization	11
2.7 Classification	11
2.8 Decision Tree (DT)	12
2.9 Decision Tree for Classification	12
2.10 Decision Tree algorithms	13
2.10.1 ID3 Algorithm	13
2.10.2 C4.5 Algorithm	14
2.11 Impurity Measure	14
2.11.1 Entropy	14
2.11.2 Information Gain	15
2.11.3 Gain Ratio	15
2.12 Parameters of Decision Tree Algorithm	16
2.13 Random Forest (RF)	16
2.14 Random Forest for Classification	17
2.15 Parameters of Random Forest Algorithm	18
2.16 Artificial Neural Networks (ANN)	19
2.17 Types of Activation Function	20
2.18 Architecture of Neural Networks	22
2.18.1 Feed Forward Neural Networks	22
2.18.2 Back Propagation Neural Networks	23
2.19 Parameters of Artificial Neural Networks Algorithm	25
2.20 Support Vector Machine (SVM)	26
2.21 Support Vector Machine for Classification	27
2.22 The Kernel Function	27
2.23 Types of Support Vector Machine	29

2.23.1 Linear SVM	29
2.23.2 Non- Linear SVM	30
2.24 Parameters of Support Vector Machine Algorithm	30
2.25 k-Nearest Neighbor Algorithm (k-NN)	31
2.26 k-Nearest Neighbor for Classification	31
2.27 Ball-Tree Algorithm	33
2.28 Distance Function	33
2.29 Parameters of k-nearest Neighbor Algorithm	34
2.30 Performance of Evaluation Criteria	34
2.30.1 Accuracy	35
2.30.2 Precision	35
2.30.3 Recall	35
2.30.4 F1-Score	35
<b><i>Chapter Three: The Proposed model</i></b>	
3.1 Introduction	36
3.2 Architecture of proposed Model	36
3.3 Data pre-processing stage	37
3.3.1 Handle Miss Values	38
3.3.2 Label Encoding	38
3.4 Split the dataset	38
3.5 Data Normalization	38
3.6 Classification Using data mining Techniques	39
3.6.1 Decision Tree	39
3.6.2 Random forest	40
3.6.3 Artificial Neural Network	41
3.6.4 Support Vector Machine	43
3.6.5 K-Nearest Neighbour	44
<b><i>Chapter Four: Experimental results and Evaluation</i></b>	
4.1 Introduction	46
4.2 Data sets	46
4.2.1 Iraqi liver patient dataset	46
4.2.2 Indian Liver Patient Dataset (ILPD)	48
4.3 Proposed Model Implementation	51
4.3.1 Data Pre-processing result	51
4.3.2 Split the Dataset	54
4.3.3 Data Normalization	54
4.3.4 Classification Model Results	56
4.3.4.1 The First Classifier Results	56
4.3.4.2 Evaluation of the First Classifier	61
4.3.4.3 The Second Classifier Results	62
4.3.4.4 Evaluation of the Second Classifier	69

4.3.4.5 The Third Classifier Results	70
4.3.4.6 Evaluation of the Third Classifier	72
4.3.4.7 The fourth Classifier Results	73
4.3.4.8 Evaluation of the Fourth Classifier	74
4.3.4.9 The Fifth Classifier Results	75
4.3.4.10 Evaluation of the Fifth Classifier	77
4.4 proposed algorithm vs. related works	78
<b><i>Chapter five: conclusion and future work</i></b>	
5.1 introduction	80
5.2 conclusions	80
5.3 suggestions for future works	81

### List of Abbreviations

<b>Abbreviations</b>	<b>Description</b>
ANN	Artificial Neural Network
AST	Aspartate aminotransferase
ALT	Alanine transaminase
ALP	Alkaline phosphatase
CLS	Conceptual Learning System
DT	Decision Tree
ID3	Iterative DiChaudomiser3
ILPD	Indian Liver Patient Dataset
K-NN	K-Nearest Neighbor
KDD	Knowledge Discovery Databases
RF	Random Forest
SVM	Support Vector Machine

### List of Figures

<b>Figure No.</b>	<b>Figure Title</b>	<b>Page No.</b>
2.1	Liver disease stages	7
2.2	Steps in KDD Process	9
2.3	Decision Tree Classification	13
2.4	The General Illustration of Random Forest	17
2.5	the Most Famous Used Activation Functions	21
2.6	fully connected Multi-layer perceptron	23



2.7	Neural Network with three layers	23
2.8	Hyperplanes example of an SVM Algorithm	27
2.9	Kernel function mapping of SVM Algorithm	28
2.10	Classify unknown sample according for its neighbour	32
2.11	An example of the ball-tree structure	33
3.1	Block Diagram of the Proposed model	37
4.1	The ratio of liver function test results for Iraqi liver patient dataset	47
4.2	The ratio of Gender of liver patients for Iraqi liver patient dataset	47
4.3	The ratio of liver function test results for ILPD	49
4.4	The ratio of Gender of liver patients for ILPD	50
4.5	Attribute before and after handle missing value for ILPD	52
4.6	Results of Z- Score normalization for Iraqi liver patient data after apply StandarScaler tool.	55
4.7	Results of Z- Score normalization for Indian liver patient data after apply StandarScaler tool.	56
4.8	The rules of Decision Tree for Iraqi liver patient dataset	57
4.9	The diagram of Decision Tree for Iraqi liver patient dataset	58
4.10	Result of the testing phase	59
4.11	The rules of Decision Tree for (ILPD)	59

4.12	The diagram of the Decision Tree for (ILPD)	60
4.13	Result of the testing phase for (ILPD)	61
4.14	The values Confusion Matrix of DT algorithm for Iraqi liver patient dataset	61
4.15	The values Confusion Matrix of DT algorithm for Indian Liver Patient Dataset	61
4.16	The evaluation criteria of the first classifier for two liver dataset	62
4.17	The rules of the first tree in Random forest for iraqi liver patient dataset	64
4.18	The diagram of the first tree in Random forest for Iraqi liver patient dataset	64
4.19	The rules of the last tree (tree 89) in Random forest for Iraqi liver patient dataset	65
4.20	The diagram of the last tree (tree 89) in Random forest for Iraqi liver patient dataset	66
4.21	the result test samples of the testing phase	66
4.22	The rules of the tree 13 in Random forest for (ILPD)	67
4.23	The diagram of the tree 13 in Random forest for (ILPD)	68
4.24	The rules of the tree 20 in Random forest for (ILPD)	69
4.25	The diagram of the tree 20 in Random forest for (ILPD)	69
4.26	The result test samples of the testing phase for (ILPD).	69
4.27	The values Confusion Matrix of RF algorithm for Iraqi liver patient dataset	70

4.28	The values Confusion Matrix of RF algorithm for Indian Liver Patient Dataset	70
4.29	The evaluation criteria of the second classifier for two liver dataset	70
4.30	The result test samples of the testing phase	71
4.31	The result test samples of the testing phase.	72
4.32	The values Confusion Matrix of ANN algorithm for Iraqi liver patient dataset	72
4.33	The values Confusion Matrix of ANN algorithm for Indian liver patient dataset	72
4.34	The evaluation criteria of the third classifier for two liver dataset	73
4.35	The result test samples of the testing phase	74
4.36	The result test samples of the testing phase.	74
4.37	The values Confusion Matrix of SVM algorithm for Iraqi liver patient dataset	74
4.38	The values Confusion Matrix of SVM algorithm for Indian liver patient dataset	75
4.39	The evaluation criteria of the fourth classifier for two liver dataset	75
4.40	The result test samples of the testing phase for the Iraqi liver patient dataset	76
4.41	The result test samples of the testing phase for the Indian Liver Patient Dataset	76
4.42	The values of Confusion Matrix KNN algorithm for Iraqi liver patient dataset	77
4.43	The values of Confusion Matrix of KNN algorithm for Indian liver patient dataset	77

4.45	The evaluation criteria of the fifth classifier for two liver dataset	78
------	--	----

### List of Tables

<b>Table No.</b>	<b>Table Title</b>	<b>Page No.</b>
2.1	Confusion Matrix	34
4.1	Attributes of Iraqi liver patient dataset	46
4.2	Some samples of Original Iraqi liver patient data	48
4.3	Attributes of Indian Liver Patient Dataset (ILPD)	49
4.4	Some samples of Original Indian Liver Patient Dataset (ILPD)	50
4.5	Samples of the original Iraqi liver patient data before and after convert labels.	53
4.6	Samples of the original Indian Liver Patient Dataset before and after convert labels.	53
4.7	The result of splitting for two dataset	54
4.8	Description of an Iraqi liver patient dataset	54
4.9	Description of an Indian Liver Patient Dataset	55
4.10	Parameters that are used in the building of Decision Tree for Iraqi liver patient dataset	56
4.11	Parameters that are used in the build of Random Forest for Iraqi liver patient dataset	62
4.12	Best parameters for learning ANN algorithm for Iraqi liver patient dataset	71
4.13	Best parameters for learning SVM algorithm for Iraqi liver patient dataset.	73
4.14	Best parameters for testing the KNN algorithm for two dataset	76

4.15	Comparison between Other Existing Works and the Proposed Work	78
------	--	----

### **List of Algorithms**

<b>Algorithm No.</b>	<b>Algorithm Title</b>	<b>Page No.</b>
3.1	Z- Score Normalization	38
3.2	Decision Tree	39
3.3	Random Forest	40
3.4	Artificial Neural Networks	42
3.5	Support Vector Machine	43
3.6	K-Nearest Neighbor	45

# **Chapter One**

## **General Introduction**

## **Chapter one**

### **General Introduction**

#### **1.1 Introduction**

The Healthcare industry contains big and complex data that may be required in order to discover fascinating pattern of diseases and makes effective decisions with the help of different machine learning techniques. data mining techniques are used to discover knowledge in database and for medical research [1].

Data mining techniques such as classification and prediction, clustering, association rule mining and various mining methods can be useful to apply on medical data.

liver has a vital role in the human body functions from protein production to removing toxins from the body and it is essential for the survival. The failure of liver functioning leads to the death. The functioning of the liver is examined by two types of tests such as imaging test and liver function tests which help to diagnose liver diseases. Liver diseases are caused by many factors such as stress, food habits, consumption of alcohol drug intake, etc. In recent days, it could be found that it is very difficult to detect at an early stages symptoms are very hard to identify. The physician often slips to detect the liver disease which leads to improper medical treatment. Various data mining algorithms can be used to classify the various disease stages including early stage so that it could be help the physician to give the proper treatment [2].

Classification can be considered as one of the most common operations in data mining. Classification is a process that divides the dataset into specified sections and then classifies the data which is a two-phase process:

In the first phase, it develops a model based on educational datasets of databases and then creates educational dataset including records, samples, examples and things with a collection of attributes and aspects. Each sample has a specific class label [3].

In the second phase, the developed model in the previous phase is used to classify new samples [3].

Classification algorithms can be either supervised or unsupervised based on the learning mechanism. Supervised learning is implemented by set of labels defined prior in the training set. The function is mapped for new unseen data to predict the labels. Few examples are Artificial Neural Network, Bagging, Boosting, Naive Bayes, Kernel-based classifiers, Nearest Neighbor algorithm, Decision Trees, Random Forest, and other ensemble of classifiers. Whereas unsupervised learning identifies the missing or hidden patterns in unlabelled data without any labels. They are commonly used for dimensionality reduction of feature space. The unsupervised ensembles include clustering approaches, self-organization maps, hidden Markov models and adaptive resonance theory [4].

In this thesis, a liver disease classification system is proposed, which uses the liver function test of a possible patient along with his/her private information. Each attribute of liver function test is normalized its values using z-score normalization. The Decision Tree (DT), Random Forest (RF), Artificial Neural Network (ANN), Support Vector Machine (SVM), and K-Nearest Neighbor (KNN) algorithms is used as a classifier to classify if the person have liver disease or not

## 1.2 Related Works

This section reviews some of previous studies and explains the different techniques that are used for classify liver disease.

- **L. Alice Auxilia (2018) [5]:** proposed a paper based on different classification techniques such as Decision Tree, Naïve Bayes, Random Forest, Support Vector Machine and Artificial Neural Network. The dataset are taken from the UCI vault. It has been seen that decision tree gives better outcomes contrast with other order classification algorithms. The accuracy obtained of applying the Decision Tree Random Forest, Support Vector Machine and Artificial Neural Network was 81%, 77%, 77%, and 71%.
- **Nazmun Nahar & Ferdous Ara (2018) [6]:** proposed a paper for prediction liver disease in early stages by using different decision tree techniques J48, LMT, Random Forest, Random tree, REPTree, Decision Stump & Hoeffding Tree. The data are collected from UCI Machine Learning Repository. The analysis proves that Decision Stump provides the highest accuracy than other



techniques. Decision Stump outperforms well than other algorithms and its achieved accuracy is 70.67% while the accuracy with Random Forest is 69.30%.

- **Thirunavukkarasu K et al. (2018) [7]:** proposed a paper to predict liver disease using different classification algorithms. The algorithms used for this purpose of work is Logistic Regression, K-Nearest Neighbour and Support Vector Machines. The dataset is taken is from the Indian Liver Patient Dataset (ILPD). This is downloaded from UCL Machine Learning Repository. Accuracy score and confusion matrix is used to compare this classification algorithm. Logistic Regression and K-Nearest Neighbour have the highest accuracy with 73.97% while the accuracy with Support Vector Machines is 71.97%.
- **Nazim Razali et al. (2019) [8]:** proposed paper to predict liver disease using different classification algorithms. This study used classification and regression as data mining tasks. In classification, models are evaluated using Bayes point machine and neural network as algorithms, while regression using linear regression and Poisson regression. This study used dataset from the University of California Irvine (UCI) repository. The results showed that Bayes Point Machines algorithm is the best algorithm used for solve the problem relating to the liver disease with accuracy 70% while neural network achieved accuracy with 66%.
- **A.K.M Sazzadur Rahman et al. (2019) [9]:** proposed paper to evaluate the performance of different Machine Learning algorithms in order to reduce the high cost of chronic liver disease diagnosis by prediction. They collected a dataset from the UCI Machine Learning Repository. They used six algorithms Logistic Regression, K Nearest Neighbors, Decision Tree, Support Vector Machine, Naïve Bayes, and Random Forest. The performance of different classification techniques was evaluated and The accuracy that obtained was 75%, 74%, 69%, 64%, 62% and 53% for LR, RF, DT, SVM, KNN and NB.
- **Md. Shafiu Azam et al. (2020) [10]:** proposed paper to predict liver disease using some efficient classification algorithms: Random Forest, Perceptron, Decision Tree, K-Nearest Neighbors (KNN), and Support Vector Machine (SVM). These classification

algorithms are applied to the original liver patient datasets collected from the UCI repository, then analysed features and tweaked to improve the performance of this predictor and made a comparative analysis among the classifiers. The accuracy before Apply feature selection was 60%, 39%, 64%, 66%, and 71% for DT, P, RF, KNN, and SVM. The accuracy after apply feature selection became 72%, 66%, 73%, 74%, and 72%.

- **Rashid Naseem et al. (2020) [11]:** In this work, Researchers are trying to project a model for the early prediction of liver disease utilizing various machine learning approaches. They used ten classifiers including A1DE, NB, MLP, SVM, KNN, CHIRP, CDT, Forest-PA, J48, and RF to find the optimal solution for early and accurate prediction of liver disease. The datasets utilized in this study are taken from the UCI ML repository and the GitHub (ILPD) repository. It has been seen that the RF classifier has highest accuracy for dataset from UCI ML repository and accuracy was 72.17%,58.26% and 62,89% for RF, SVM, KNN. SVM has highest accuracy for Dataset from GitHub repository (ILPD) and accuracy was 71.35%, 69.29%, and 64.15% for SVM, RF, and KNN.
- **Jagdeep Singh et al. (2020) [12]:** proposed paper to predict liver disease using classification and feature selection technique. The implementation of proposed work is done using (ILPD) using different classification algorithms such as such as Random Forest, K-Nearest Neighbor, SMO, Logistic Regression, Naïve Bayes, and J48. The accuracy was achieved by RF, and KNN before using feature selection technique 71.53%, and 64.15%, while the accuracy after apply feature selection technique 71.87%, and 67.41%.

### 1.3 Problem Statement

Patients with the liver disease continue to increase and the symptoms of the disease are difficult to detect. Therefore many people suffer from liver damage but they feel healthy, it causes many medical practitioners to often fail to detect the disease. Therefore accurate detection is necessary to help the medical practitioner to give proper medication and medical treatment. So the medical classification model will help a doctor to automatic classification and reduces the workload.

## **1.4 Aim of Thesis**

This thesis aims to design and implement a liver disease classification model able to accurately classify if the person have liver disease or not based on the 10 important attributes of liver disease using a DT, RF, ANN, SVM, and KNN.

## **1.5 Outline of Thesis**

The rest chapters in this thesis are organized as follows:

### **Chapter Two: Theoretical Background**

This chapter clarifies the definition and Knowledge discovery in databases (KDD) steps, it also explains the Data mining techniques used for classification.

### **Chapter Three: The Proposed System**

This chapter describes the proposed classification model with its design and implementation.

### **Chapter Four: Experimental Results and Evaluation**

This chapter shows the implementation results of the proposed model steps and evaluates these results.

### **Chapter Five: Conclusions and Suggestions for Future work**

This chapter presents the conclusions of this work. Furthermore, it provides suggestions for future work.