

# Study on Classification Algorithms for Multi-relational Data Mining

Mahfuza Mallika

Lecturer in CSE, Centre for General Education, Bangladesh Islami University, Dhaka, Bangladesh

DOI: <https://doi.org/10.51244/IJRSI.2025.1210000276>

Received: 30 October 2025; Accepted: 06 November 2025; Published: 19 November 2025

## ABSTRACT

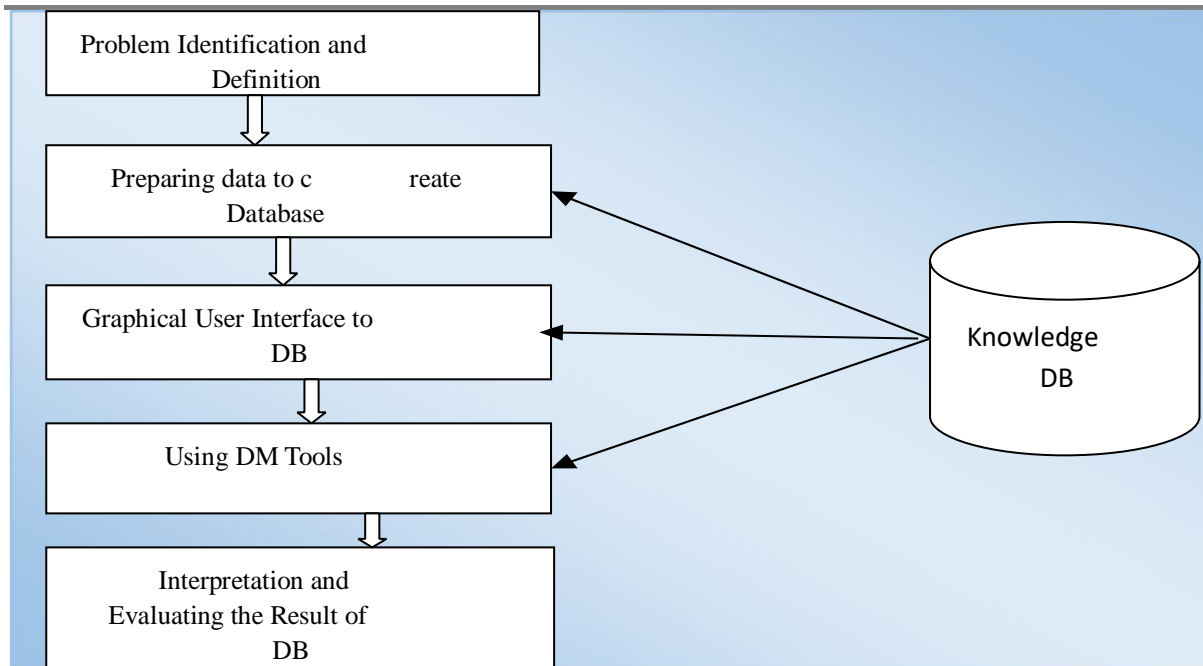
Data Mining (DM) constitutes a fundamental stage in the Knowledge Discovery in Databases (KDD) process, emphasizing the systematic analysis of large-scale datasets to uncover meaningful patterns, trends, and relationships. As data continues to grow in volume and complexity, the application of advanced analytical techniques has become essential for transforming raw data into actionable knowledge. DM employs a range of methods to address analytical challenges, including classification algorithms, association rule mining, and neural network approaches. This study investigates the effectiveness of various classification algorithms namely: ID3, C4.5, J48 and the general Decision Tree methodology in solving classification problems within data mining tasks. A benchmark dataset from the UCI Machine Learning Repository is utilized to illustrate the practical application of these algorithms. The WEKA software tool is employed for data preprocessing, model development, and performance evaluation through metrics such as accuracy and predictive power. The experimental results highlight the capability of classification techniques to categorize data points efficiently and extract valuable insights. Overall, the study underscores the critical role of classification-based data mining techniques in enhancing knowledge discovery and supporting informed decision-making across diverse domains.

**Keywords:** Classification and Classification Algorithm, Data Mining Techniques, Decision Tree, Weka, Application of J48 algorithm, ID3 Algorithm, C4.5 Algorithm, MLP, DRDM, UCI.

## INTRODUCTION

Data Mining is one of the most popular topics today and is considered an emerging research area across many disciplines. It plays a crucial role in efficient data analysis and provides accurate reports for decision-making. Unlike rigid rules, data mining applies flexible methods to analyze datasets. In simple terms, it is the process of extracting useful information from large volumes of data by combining techniques from statistics, artificial intelligence, neural networks, and related fields. Currently, data mining is widely applied in diverse domains such as business, science, engineering, and medicine. Some common applications include mining credit card transactions, analyzing stock market movements, ensuring national security, and evaluating clinical trials. Retail and financial companies also rely on data mining to recognize customer trends, forecast stock prices, and analyze interest rates influenced by government policies.

Data mining uses various algorithms and performs multiple tasks depending on the problem at hand, with the goal of predicting accurate results. These tasks involve identifying patterns within large datasets and can be categorized into classification, regression, clustering, summarization, and association rules, among others. Overall, the data mining process is highly valuable for analyzing datasets and uncovering meaningful insights. The general phase of the mining process is illustrated in Figure 1.



**Figure 1: The process of Data Mining**

## 1.2 Objectives of this Study

The primary objectives of this study are to investigate the effectiveness of various classification algorithms, with a particular focus on decision tree methods, in analyzing large datasets. The research aims to evaluate the role of decision tree algorithms, specifically ID3 and C4.5, in generating accurate and efficient classification rules. Additionally, it explores the application of multi-relational classification techniques using decision approaches to handle complex relational data. The study also analyzes the use of WEKA as a practical data mining tool for performing classification tasks and emphasizes its role in facilitating experimental analysis. Furthermore, it examines the significance of association rule mining in uncovering hidden relationships within relational databases, contributing to the broader field of knowledge discovery.

## METHODOLOGY

This study adopts a quantitative and analysis research methodology, employing a systematic framework to evaluate and compare classification algorithms for multi-relational data mining. The research follows the Knowledge Discovery in Databases (KDD) process, which integrates Data Mining (DM) techniques to analyze large datasets and uncover meaningful patterns and relationships. A dataset from the UCI Machine Learning Repository is selected as the primary data source due to its reliability, accessibility, and relevance to classification tasks. Prior to analysis, the dataset undergoes pre-processing using the WEKA software environment. This stage includes handling missing values, normalizing data, and performing attribute selection to enhance data quality, consistency, and suitability for mining tasks. Four classification algorithms, namely ID3, C4.5, J48, and a general Decision Tree approach, are implemented within WEKA to construct predictive models. Each algorithm is applied to the pre-processed dataset using WEKA's built-in modules, ensuring methodological consistency and reliability. To assess model performance, a 10-fold cross-validation technique is employed, providing a robust and unbiased estimate of each algorithm's predictive capability. Model performance is evaluated using standard classification metrics, including True Positive (TP) rate, False Positive (FP) rate, Precision, Recall, F-measure, and the confusion matrix, which collectively assess the accuracy, reliability, and efficiency of each algorithm. A comparative analysis is then conducted to determine the most effective algorithm for multi-relational data mining tasks. The findings are expected to contribute to a deeper understanding of how DM techniques can be effectively applied for knowledge extraction and decision-making in complex, multi-relational datasets.

## 1.4 Related Work

Classification is a fundamental task in data mining that categorizes large datasets into predefined classes using algorithms such as decision trees, neural networks, and rule-based methods (Han, Pei, and Kamber, 2011). Among these, decision tree algorithms have gained significant popularity due to their simplicity, interpretability, and the ease with which rules can be extracted (Quinlan, 1993). Decision trees classify relational data based on decision rules, and well-known algorithms such as ID3 and C4.5 are widely applied for tree construction (Quinlan, 1986; Quinlan, 1993). Recent studies have extended them to multi-relational classification, enabling analysis across interconnected relational tables (Knobbe, Blockeel, and Siebes, 1999).

For practical implementation, the WEKA software provides a robust environment for classification and association analysis, including the J48 algorithm (an implementation of C4.5) (Witten, Frank, and Hall, 2011). In addition to classification, association rule mining is an important data mining technique for discovering hidden relationships in large relational or multi-relational databases. This method, first introduced by Agrawal, Imelinski, and Swami (1993), identifies correlations among seemingly unrelated attributes and has proven especially valuable in domains such as market basket analysis. Together, decision tree classification and association rule mining provide complementary approaches for predictive modeling and pattern discovery.

## 2. Definition of Classification and Classification Algorithms

Classification is one of the fundamental techniques in data mining used to analyze and interpret large-scale data repositories. It involves assigning objects or instances to predefined categories based on their similarity to existing data patterns and is typically performed using a model trained from historical data (Han, Pei, & Kamber, 2011). A classifier is essentially this model, which predicts categorical labels, such as determining whether a financial transaction is genuine or fraudulent, classifying a bank loan application as safe or risky, predicting customer responses to marketing campaigns as yes or no, or diagnosing medical causes such as blood group classification (A, B, O) or engine faults based on symptoms (Kantardzic, 2019).

The process of classification typically follows two steps. In the training phase (Supervised Learning), a classification algorithm, commonly a decision tree, neural network, or rule-based approach, analyzes labeled data to generate classification rules and models. In the testing phase (unsupervised evaluation), unseen data (test data) is analyzed using these rules to validate prediction accuracy (Tan, Steinbach, and Kumer, 2019).

A classification algorithm can be understood as a procedure for selecting the most suitable hypothesis from a set of alternatives that best explain observed data. Among the most influential algorithms are ID3 (Iterative Dichotomiser3) and C4.5, introduced by Quinlan (1986, 1993), which construct decision tree models from training data. Each case is represented by a set of attributes and their values, and the classifier model predicts the most appropriate class for new instances. The ultimate goal of classification is accurate decision-making and prediction, which is central to many fields ranging from finance to healthcare.

### 2.1 ID3 Algorithm

The ID3 (Iterative Dichotomiser 3) algorithm is one of the most widely used methods for constructing decision trees in data mining and machine learning. The algorithm begins with a set  $S$ , which contains a list of attributes relevant to the classification process. To determine the most effective way to split the data, ID3 evaluates each attribute by calculating entropy and information gain. The attribute with the maximum information gain (or equivalently the maximum entropy) is selected as the splitting criterion. Once the attribute is chosen, the dataset is partitioned into subsets such that each subset contains instances with similar attribute values. Each subset forms a node at the decision tree, eventually leading to classification at the class level. The process is applied recursively until all instances are classified or no further attributes remain. The final decision tree consists of node and class levels that represent the classification rules derived from the dataset. By searching through the subsets and applying attributes step by step, ID3 effectively generates a decision structure that can be used for prediction and classification task (Quinlan 1986).

## 2.2 C4.5 Algorithm

The C4.5 algorithm, developed as an extension of the ID3 algorithm, is widely applied for constructing decision trees from training datasets. Similar to ID3, it classifies a dataset  $S$  into subsets, denoted as  $S = \{s_1, s_2, \dots, s_n\}$ , where each subset belongs to the same class and represents a group of samples with defined attribute values or features. These subsets are then used to progressively build the decision tree.

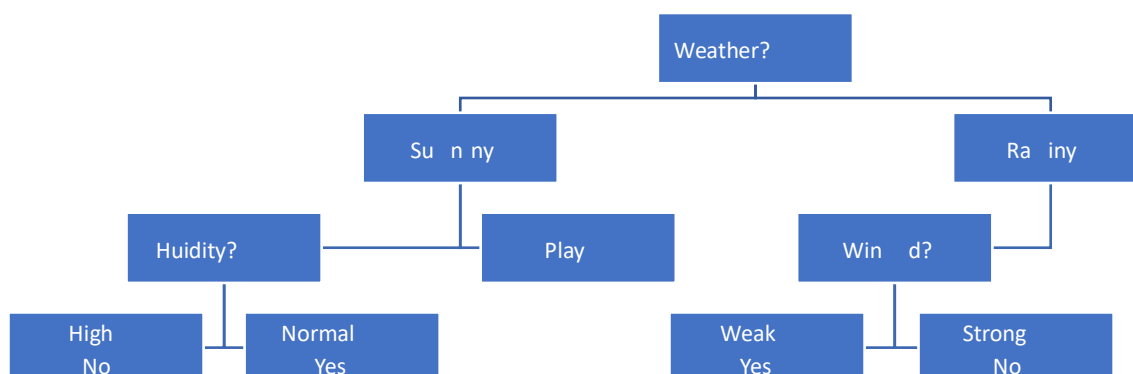
A key improvement in C4.5 over ID3 is the use of normalized information gain (gain ratio) to select the splitting attribute, ensuring a more balanced and accurate decision-making process. The attribute that provides the highest normalized information gain is chosen as the root or internal node, and this recursive process continues until all samples are classified. The resulting decision tree thus represents a hierarchical model of decision rules derived from the training data, making C4.5 an effective algorithm for classification tasks (Quinlan, 1993).

In addition, C4.5 can handle both continuous and discrete attributes, automatically determining threshold values for continuous data splits. It also incorporates pruning techniques to reduce overfitting, ensuring that the final decision tree generalizes better to unseen data. Furthermore, C4.5 can effectively manage datasets with missing attribute values, which enhances its usability in real-world applications where data may be incomplete (Quinlan, 1993).

## 2.3 Decision Tree

Decision trees are widely used models in data mining, statistics, and machine learning for representing a set of rules that guide the process of decision-making and classification. They serve as the foundation for several algorithms such as J48, ID3, C4.5, and CART, which have been extensively applied in classification and prediction tasks. A decision tree is structured with a root node, branches, and levels. The root node represents the initial test attribute, from which branches emerge. Each branch denotes a conjunction of feature values that eventually leads to a terminal node or a leaf, representing the class label or outcome. Depending on the type of target variable, decision trees are categorized into classification trees and regression trees. A classification tree is used when the target variable contains categorical values, whereas regression trees are employed when the target variable consists of continuous values. In terms of structure, decision trees may be binary (each node splits into two branches) or multi-way (nodes split into more than two branches).

A decision tree is composed of internal and external nodes connected by arcs. Each internal node functions as a decision-making unit, where input attributes are tested to determine the next child node to be visited. The arcs represent the possible values of these attributes. External nodes or leaves contain no further splits and are associated with a specific class label or a probabilistic distribution of classes. The hierarchical structure of nodes and branches makes decision trees interpretable and effective, contributing to their broad application in predictive analytics and knowledge discovery (Quinlan, 1993; Breiman et al., 1984). For example,



## Figure 2: Decision Tree

**Root Node:** Weather; **Branches:** Conditions like Sunny, Rainy; **Leaves:** Yes/No (class labels); **Internal Nodes:** Weather, Humidity, Wind; **External Nodes:** Yes/No (final decision)

### 3.1 Training Dataset

We use data on hepatitis from the UCI Machine Learning Repository for this study. The dataset details are given below:

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances</b>	155	<b>Area:</b>	Life
<b>Attribute Characteristics:</b>	Categorical, Integer, Real	<b>Number of Attributes:</b>	19	<b>Date Donated</b>	1988-11-01
<b>Associated Tasks:</b>	Classification	<b>Missing Values</b>	Yes	<b>Number of Web Hits:</b>	237466

**Table 1: Details of Hepatitis Datasets**

Hepatitis is a medical condition characterized by the inflammation of liver tissue, which can ultimately lead to severe liver damage. There are several types of hepatitis, including Hepatitis A, B, C, D, and others\*\*. These types are primarily caused by different viruses, although additional factors such as certain drugs, toxins, and excessive alcohol consumption can also contribute to the development of the disease (World Health Organization, 2023, Center for Disease Control and Prevention [CDC], 2022).

For experimental and analytical studies, a Hepatitis dataset is employed to classify data and generate predictive outcomes using machine learning and data mining algorithms. This dataset contains 20 attributes and 155 instances, providing comprehensive information relevant to the classification of hepatitis conditions. Among these attributes, Bilirubin is defined as a continuous attribute, which plays an essential role in evaluating liver function and distinguishing between different classes of data (UCI Machine Learning Repository, n.d.).

### 3.2 WEKA(Waikato Environment for knowledge Analysis)

WEKA (Waikato Environment for Knowledge Analysis) is an open-source machine learning software with a vast collection of algorithms for data mining. It was developed by the University of Waikato in New Zealand, and it's written in Java. It supports different data mining tasks, such as pre-processing, classification, regression, clustering, and visualization in a graphical user interface that makes it easy to use. For each of these tasks, WEKA provides built-in machine learning algorithms, which allow users to quickly test their ideas and deploy models without writing any code. Users need to have knowledge of different algorithms so that they can choose the right one for their task. The benefits include free availability under the GNU (General Public License), a comprehensive collection of algorithms, use of the Java Programming Language, support for modern computing platforms, the ability to compare different approaches, and more (Singhal, S. et al., 2013).

### 3.3 Application of J48 Algorithm in WEKA for Data Classification

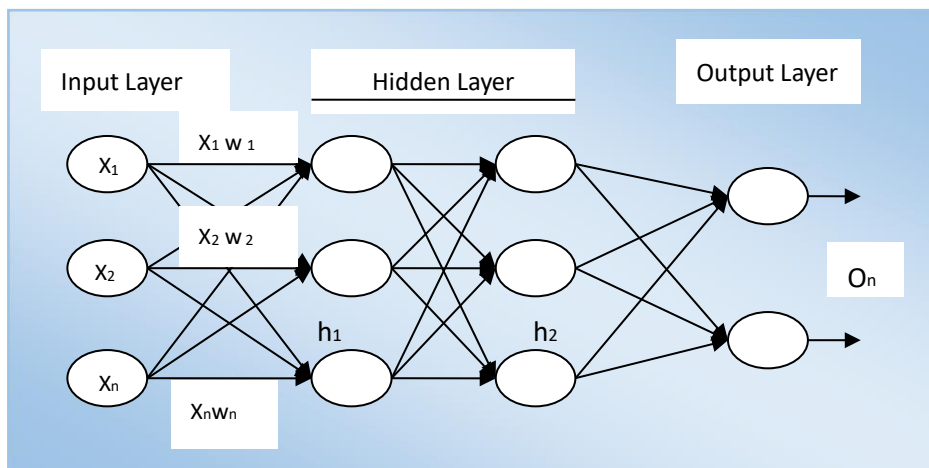
WEKA is a widely used open-source data mining software tool that supports various machine learning algorithms for data analysis and predictive modeling. Among these, the J48 algorithm, developed by Ross Quinlan, is one of the most significant and frequently applied. J48 is an extension of the ID3 (Iterative Dichotomiser 3) algorithm and is designed to generate efficient decision trees for both small and large datasets.

The decision tree model constructed by J48 offers a highly interpretable and convenient structure for classifying complex and critical datasets. By utilizing the decision tree, the algorithm enables accurate prediction of target variables based on input attributes.

One of the notable strengths of J48 lies in its advantageous features, which include handling of missing values, decision tree pruning to prevent overfitting, management of continuous attribute value ranges, and derivation of classification rules for enhanced interpretability (Quinlan, 1993). The algorithm operates recursively, classifying each level of the dataset as precisely as possible, ultimately leading to a robust classification model.

### 3.4 Multilayer perception algorithm

The Multilayer Perceptron (MLP) algorithm is applied to construct an artificial neural network that is capable of generating complex outputs from simple processing mechanisms. It operates as a feed-forward network and is composed of three primary layers: the input layer, one or more hidden layers, and the output layer. The input layer receives the external data, where each input is associated with a corresponding weight that determines its level of influence and accelerates the activation process. These weighted inputs are transmitted simultaneously to the hidden layer(s), where they are processed through activation functions to extract non-linear patterns and relationships. The outputs generated from one hidden layer are then passed forward as inputs to the subsequent hidden layer(s), allowing the network to learn and represent more complex functions. Finally, the output layer produces the final response of the network. This layered structure enables the MLP to approximate complex functions and is widely utilized in various machine learning and pattern recognition tasks (Haykin, 1999; Bishop, 2006). For example,



**Figure 3: A multi-layer feed forward neural network**

The inputs are  $x_1, x_2, x_3$  and the weights are  $w_1, w_2, w_3$ . The net input (weighted sum) is:

$$f(x) = \sum x_i w_i.$$

The back-propagation algorithm is used in MLP to adjust the weights and enable learning. The back-propagation algorithm performs learning on a multilayer feed-forward neural network. At its core, back propagation is simply an efficient and exact method for calculating all the derivatives of a single target quantity (such as pattern classification error) with respect to a large set of input quantities (such as the parameters or weights in classification rules)

### 4. The Output of Classification using Weka tools

The J48 algorithm was applied to the Hepatitis dataset obtained from the UCI Machine Learning Repository

(Dua & Graff, 2019). Initially, the algorithm utilized the training dataset to construct a decision tree model. Once the model was developed, it was subsequently applied to the dataset objects to classify and analyze the instances effectively. This process demonstrates how the J48 algorithm can transform raw data into a structured decisionmaking model, enabling accurate prediction and classification within medical datasets.

Consider the J48 algorithm running on the hepatitis dataset in WEKA. For this dataset, since we get two classes, we have a 2 times 2 confusion matrix. There were 155 classified instances in total, with 99 instances correctly classified and 56 instances incorrectly classified. The confusion matrix is telling the following:

1. The decision tree correctly classified 65 instances as **Class A** and misclassified 20 instances as **Class B**.
2. The decision tree correctly classified 34 instances as **Class B** and misclassified 36 instances as **Class A**.

The confusin matrix of J48 algorithm are given below:

====		Confusion Matrix		====	
A	B	<	--	---	Classified as
65	20				A = No
36	34				B = Yes

**Figure 4: Classifier output**

The **True Positive (TP) Rate**: In the confusion matrix, this is the diagonal element divided by the sum over the **relevant row**, i.e., 65/85 for Class 'No' and 34/70 for Class 'Yes'. The **False Positive (FP) Rate**: In the confusion matrix, this is the column sum of class x minus the diagonal element, divided by the row sums of all other classes; i.e., 36/70 for the 'No' class and 20/85 for the 'Yes' class. The **Precision** : In this confusion matrix, this ratio is the diagonal element divided by the sum over the relevant column; i.e., 65/101 for the 'No' class and 20/54 for the 'Yes' class. **Recall**: It is the same as True Positive (TP) Rate. The **F-Measure**: It is simply  $2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$ . By these measures, we can compare the classifiers and get the classifiers' predictions. The confusion matrix table for accuracy is given below:

Actual Class	Predicted Class			
		A	B	Total
	A	65	20	85
	B	36	34	70
	Total			155

**Table 2: The confusion matrix table for accuracy**

Accuracy =  $(TP_A + TP_B) / (\text{Total number of classification})$ , i.e. Accuracy =  $(65 + 34) / 155 = 0.639$ .

## 5. Multi-Relational Data Mining (MRDM)

Data mining algorithms are traditionally applied to a single data table to extract patterns, relationships, or knowledge. In contrast, Multi-Relational Data Mining (MRDM) extends these techniques to operate across

multiple interrelated tables within a relational database. MRDM represents an advanced form of data mining that leverages the rich information embedded in multiple relations to derive more comprehensive and insightful results. MRDM has applications across various domains, including sports analytics (e.g., IPL), Knowledge Discovery in Databases (KDD), statistics, and machine learning, offering optimal solutions by integrating information from multiple sources (Pal, S. K., & Mitra, S., 1992; Džeroski, S. 2003).

In a relational database, relations (tables) are connected via keys, which establish logical links between them. For instance, in a database containing patient information, the Patient and Hepatitis tables may be interconnected, with each table defined both extensionally (through stored data) and intentionally (through logical rules).

MRDM focuses on discovering multi-relational patterns, which include: Multi-relational classification rules (rules that classify entities based on attributes spread across multiple tables), Multi-relational regression rules (rules that predict continuous values using data from several interconnected relations), and Multi-relational association rules (rules that identify associations among entities spanning multiple tables). By integrating data across multiple tables, MRDM enables more sophisticated analysis than single-table approaches, uncovering patterns that would otherwise remain hidden in isolated datasets.

Multi-Relational Data Mining can be effectively applied in the field of neural network systems. The framework of multi-relational data mining is based on searching for interesting patterns within relational databases. By combining data mining methods with neural network models, the efficiency of data mining techniques can be significantly improved, and this approach has been widely adopted in various applications (Zhang et al., 2019).

## CONCLUSION

The primary objective of this study was to evaluate the performance of the J48 classification algorithm using the Weka data mining tool, with a focus on its application to the Hepatitis dataset obtained from the UCI Machine Learning Repository. Through systematic experimentation, the algorithm was trained and tested on the dataset to assess its ability to classify instances effectively and generate accurate predictive models.

The evaluation process involved the computation of several key performance metrics, including True Positive (TP) rate, False Positive (FP) rate, Precision, Recall, and F-measure. These metrics provided a comprehensive understanding of the algorithm's classification efficiency and reliability. Additionally, the confusion matrix was analyzed in detail to further validate the accuracy of the model, offering insights into how well the algorithm distinguished between different classes within the dataset. This step was crucial for identifying both the correctly classified instances and potential misclassifications.

The findings indicate that the J48 algorithm is not only efficient in handling medical datasets but also provides clear, interpretable decision tree models, which are particularly valuable in medical research and decision making processes. Its ability to manage missing values, prune decision trees, and generate understandable classification rules makes it suitable for practical applications in healthcare analytics, where transparency and interpretability are essential.

Moreover, this research highlights the potential of J48 as a baseline algorithm for medical classification tasks. Future studies could extend this work by incorporating additional classification techniques, such as Random Forest, Support Vector Machines (SVM), or ensemble methods, to compare performance and enhance predictive accuracy. Parameter tuning, feature selection techniques, and cross-validation strategies may also be applied to optimize model performance further. Additionally, expanding the dataset or applying the algorithm to other healthcare domains could provide broader insights into its generalizability and robustness.

Overall, the results of this study demonstrate that the J48 classification algorithm, when implemented through Weka, is a robust approach for analyzing complex medical datasets, effectively supporting both academic research and real-world clinical applications.

## REFERENCES

1. Agrawal, R., Imielinski, T., & Swami, A. (1993). Mining association rules between sets of items in large databases. *ACM SIGMOD Record*, 22(2), 207–216.
2. Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
3. Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees*. CRC Press.
4. Centers for Disease Control and Prevention. (2022). Hepatitis: Overview and statistics. Retrieved from <https://www.cdc.gov/hepatitis>
5. Dua, D., & Graff, C. (2019). UCI machine learning repository. University of California, Irvine. Retrieved from <https://archive.ics.uci.edu>
6. Džeroski, S. (2003). Multi-relational data mining: An introduction. *ACM SIGKDD Explorations Newsletter*, 5(1), 1–16.
7. Han, J., Kamber, M., & Pei, J. (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann.
8. Han, J., Kamber, M., & Pei, J. (2012). *Data mining: Concepts and techniques*. Elsevier.
9. Habrard, A., Bernard, M., & Jacquenet, F. (2003, October). Multi-relational data mining in medical databases. In *Conference on Artificial Intelligence in Medicine in Europe* (pp. 365–374). Springer Berlin Heidelberg.
10. Haykin, S. (1999). *Neural networks: A comprehensive foundation* (2nd ed.). Prentice Hall.
11. Kaczmarek, B. L. J., & Knobbe, A. J. (2006). *Multi-relational data mining* (Vol. 145). IOS Press.
12. Kanodia, J. (2005). Structural advances for pattern discovery in multi-relational databases.
13. Kantardzic, M. (2019). *Data mining: Concepts, models, methods, and algorithms*. Wiley.
14. Knobbe, A., Blockeel, H., & Siebes, A. (1999). Multi-relational decision tree induction. In *Proceedings of the Third European Symposium on Principles of Data Mining and Knowledge Discovery* (pp. 1–12).
15. Morariu, D., Crețulescu, R., & Breazu, M. (2017). The WEKA multilayer perceptron classifier. *International Journal of Advanced Statistics and IT&C for Economics and Life Sciences*, 7(1), 1–7.
16. Murtagh, F. (1991). Multilayer perceptrons for classification and regression. *Neurocomputing*, 2(5–6), 183–197.
17. Padhy, N., & Panigrahi, R. (2012). Multi-relational data mining approaches: A data mining technique. *arXiv preprint arXiv:1211.3871*.
18. Pal, S. K., & Mitra, S. (1992). Multilayer perceptron, fuzzy sets, and classification. *IEEE Transactions on Neural Networks*, 3(5), 683–697.
19. Quinlan, J. R. (1986). Introduction to decision trees. *Machine Learning*, 1(1), 81–106.
20. Quinlan, J. R. (1993). *C4.5: Programs for machine learning*. Morgan Kaufmann.
21. Singhal, S., & Jena, M. (2013). A study on WEKA tool for data preprocessing, classification and clustering. *International Journal of Innovative Technology and Exploring Engineering (IJITEE)*, 2(6), 250–253.
22. Tan, P. N., Steinbach, M., & Kumar, V. (2019). *Introduction to data mining* (2nd ed.). Pearson.
23. UCI Machine Learning Repository. (n.d.). Hepatitis dataset. University of California, Irvine. Retrieved from <https://archive.ics.uci.edu>
24. World Health Organization. (2023). Hepatitis. Retrieved from <https://www.who.int/health-topics/hepatitis>
25. Zhang, S., Li, X., & Yang, J. (2019). Multi-relational data mining and neural network integration for efficient knowledge discovery. *International Journal of Data Mining and Knowledge Management Process*, 9(2), 15–27.