



# A Comparative Study of Machine Learning Algorithms on Structured Data

Pratibha Kulkarni<sup>1</sup>, Ms. V. Lavanya<sup>2</sup>

<sup>1</sup>MCA Student, CMR University, Bangalore, India.

<sup>2</sup>Assistant Professor, CMR University, Bangalore, India.

**ABSTRACT:** This study presents a comprehensive comparison of the performance of various machine learning algorithms on structured datasets, evaluating their accuracy, computational efficiency, memory usage, and scalability. The analysis focuses on five prominent algorithms: Decision Trees, Support Vector Machines (SVM), Random Forests, Gradient Boosting, and Logistic Regression. These algorithms are widely used in diverse domains, including healthcare, finance, and e-commerce, where accurate predictions and efficient processing are crucial.

The study's primary objective is to inform algorithm selection for specific tasks by examining the strengths and weaknesses of each algorithm. This involves assessing their performance on various structured datasets, identifying the most suitable algorithms for different domains, and providing guidelines for hyperparameter tuning to optimize model performance. A key aspect of the study is the examination of hyperparameter tuning's influence on model performance. Hyperparameters are critical components of machine learning algorithms, and their tuning can significantly impact the accuracy and efficiency of the models. The study provides recommendations for choosing suitable algorithms for specific tasks, considering the trade-offs between accuracy, computational efficiency, memory usage, and scalability.

**Keywords:** Machine Learning, Structured Data, Decision Trees, Support Vector Machines, Random Forests, Gradient Boosting, Logistic Regression, Hyperparameter Tuning.

## 1. INTRODUCTION

The significance of machine learning in data-driven industries such as healthcare, finance, and business intelligence cannot be overstated. The ability to extract valuable insights from structured data, which is characterized by well-defined features and labels, has become a crucial aspect of decision-making processes. However, the plethora of machine learning algorithms available, each with its strengths and weaknesses, poses a significant challenge in selecting the most suitable algorithm for a specific task.

This study addresses this challenge by conducting a systematic evaluation of five prominent machine learning algorithms: Decision Trees, Support Vector Machines (SVM), Random Forests, Gradient Boosting, and Logistic Regression. These algorithms are widely used in various domains, and their performance is critical in achieving accurate predictions and efficient processing.

The evaluation is based on a range of key performance metrics, including accuracy, precision, recall, F1-score, training time, and memory usage. By examining these metrics, this study provides a comprehensive framework for algorithm selection, considering the specific requirements of various applications. The study's findings will enable practitioners to make informed decisions when selecting machine learning algorithms,

ultimately leading to improved performance and efficiency in their respective domains.

## 2. OBJECTIVES

The primary objectives of this study are multifaceted and designed to provide a comprehensive understanding of the performance of machine learning algorithms on structured datasets. The first objective is to conduct a thorough comparison of the performance of various machine learning algorithms, including Decision Trees, Support Vector Machines (SVM), Random Forests, Gradient Boosting, and Logistic Regression. This comparison will involve an examination of the algorithms' strengths and weaknesses, providing insights into their suitability for different applications.

The second objective is to evaluate the algorithms' accuracy, computational efficiency, and scalability. This evaluation will provide a comprehensive understanding of the algorithms' performance, enabling practitioners to make informed decisions about algorithm selection. The study will also investigate the trade-offs between accuracy and computational cost, which is critical in many applications where computational resources are limited.

The third objective is to examine the impact of hyperparameter tuning on model performance. Hyperparameters play a critical

role in machine learning algorithms, and their tuning can significantly impact the accuracy and efficiency of the models. The study will provide guidelines for optimizing algorithm performance through hyperparameter tuning, enabling practitioners to achieve the best possible results.

The final objective is to offer evidence-based recommendations for algorithm selection, taking into account the specific requirements of various domains. The study's findings will provide a comprehensive framework for algorithm selection, enabling practitioners to make informed decisions about the most suitable algorithm for their specific task.

### 3. LITERATURE REVIEW

The existing literature on machine learning applications to structured data is vast and diverse, with numerous studies exploring the use of individual algorithms across various fields (Hastie, Tibshirani, & Friedman, 2009; Murthy, 1998). Decision Trees, Random Forests, and Support Vector Machines (SVM) are among the most used algorithms, and their performance has been extensively evaluated in specific contexts (Breiman, 2001; Cortes & Vapnik, 1995). However, comprehensive comparisons of these algorithms, encompassing multiple performance metrics and evaluation criteria, are less common (Raschka & Mirjalili, 2017).

Structured data, characterized by well-defined features and labels, presents unique challenges and opportunities for machine learning applications (James, Witten, Hastie, & Tibshirani, 2013). The tabular format of structured data necessitates a tailored approach to model selection, taking into account the specific characteristics of the data and the requirements of the application. Previous studies have highlighted the importance of ensemble methods, such as Random Forests and Gradient Boosting, which often outperform single models in terms of accuracy (Breiman, 2001; Friedman, 2001). However, these methods can be computationally demanding, which can be a significant drawback in applications where computational resources are limited (Hastie et al., 2009).

This study aims to address the existing knowledge gap by providing a comprehensive comparison of five prominent machine learning algorithms, including Decision Trees, SVM, Random Forests, Gradient Boosting, and Logistic Regression (Raschka & Mirjalili, 2017; Pedregosa et al., 2011). The study's objectives are to evaluate the performance of these algorithms across a range of metrics, including accuracy, precision, recall, F1-score, training time, and memory usage, and to provide insights into their suitability for different applications. By examining the strengths and weaknesses of each algorithm, this study aims to provide a holistic understanding of their performance and to offer evidence-based recommendations for algorithm selection.

### 4. METHODOLOGY:

#### 4.1 Data Collection

To ensure diversity and representativeness, this study utilized datasets from three distinct domains: healthcare, finance, and e-commerce. These datasets were carefully selected to provide a robust testing ground for evaluating the performance of different machine learning algorithms. The datasets comprised

a mix of categorical and numerical features, which is typical of many real-world applications.

The use of datasets from different domains allows for a more comprehensive evaluation of the algorithms' performance, as it takes into account various data characteristics and challenges. For instance, healthcare datasets often involve complex patterns and high-dimensional data, while finance datasets may require handling of categorical variables and non-linear relationships. E-commerce datasets, on the other hand, may involve large volumes of data and require efficient processing.

#### 4.2 Data Preprocessing

The preprocessing stage is a critical component of any machine learning pipeline, as it prepares the data for analysis and ensures that the algorithms can handle the data uniformly. In this study, the preprocessing stage involved several crucial steps:

1. **Missing Value Imputation:** Missing values were imputed to prevent any bias in the results. This is particularly important in datasets where missing values are common, as it can affect the accuracy of the algorithms.
2. **Normalization:** The data was normalized using Min-Max scaling to ensure that all features were on the same scale. This is essential for algorithms that rely on distance metrics, such as k-nearest neighbours or support vector machines.
3. **Categorical Variable Encoding:** Categorical variables were encoded using one-hot encoding, enabling all algorithms to handle the data uniformly. This is particularly important for algorithms that require numerical inputs, such as neural networks.

#### 4.3 Algorithms Tested

This study evaluated the performance of five prominent machine learning algorithms, each with its strengths and weaknesses:

1. **Decision Trees:** Recursive models that partition data based on feature values, known for their simplicity but prone to overfitting.
2. **Support Vector Machines (SVM):** Finds the optimal hyperplane to separate classes, suitable for high-dimensional data but computationally intensive.
3. **Random Forests:** An ensemble of Decision Trees, reducing overfitting and generally improving accuracy.
4. **Gradient Boosting:** Sequentially builds models to minimize errors, effective for complex patterns but resource-intensive.
5. **Logistic Regression:** A linear model for binary classification, efficient for simple datasets but less capable with non-linear patterns.

#### 4.4 Evaluation Metrics

To comprehensively evaluate the performance of each algorithm, this study employed the following metrics:

1. **Accuracy:** The ratio of correct predictions to total predictions.
2. **Precision:** The ratio of true positive predictions to the total predicted positives.
3. **Recall:** The ratio of true positive predictions to the actual positive instances.
4. **F1-Score:** The harmonic mean of precision and recall.
5. **Training Time:** Time taken to train the model.
6. **Memory Usage:** The memory consumed during model training and inference.

#### 4.5 Cross-Validation

To mitigate overfitting and ensure robustness, a 10-fold cross-validation was employed. This involved training and validating each model multiple times across different data splits, providing a comprehensive understanding of the algorithms' performance.

#### 4.6 Comparative Analysis of Algorithms

The table below summarizes the performance of each algorithm across the tested datasets, providing a comprehensive comparison of their strengths and weaknesses:

**Table 1:** Performance Metrics of Machine Learning Algorithms

Algorithm	Accuracy	Precision	Recall	F1-Score	Training Time (s)
Decision Tree	0.81	0.79	0.8	0.79	0.03
Support Vector Machine	0.85	0.84	0.83	0.83	0.72
Random Forest	0.88	0.87	0.86	0.87	0.12
Gradient Boosting	0.89	0.88	0.87	0.88	0.5
Logistic Regression	0.82	0.8	0.81	0.8	0.01

## 5. RESULTS AND DISCUSSION

### 5.1 Accuracy and Precision

The study's findings indicate that Gradient Boosting and Random Forests excel in accuracy and precision, with Gradient Boosting having a slight edge. These algorithms' high performance makes them suitable for applications requiring a low false positive rate, such as medical diagnoses or financial risk assessments. Support Vector Machines (SVM) also performed well, particularly in datasets with fewer features. However, SVM's longer training time and higher memory usage make it less feasible for large-scale applications.

### 5.2 Recall and F1-Score

Gradient Boosting outperformed other algorithms in recall and F1-score, demonstrating its strength in identifying true positive cases. This is particularly valuable in fields like healthcare, where false negatives can have severe consequences. Logistic Regression, while efficient in training time, showed lower recall and F1-score, making it less ideal for applications where identifying all positive cases is crucial.

### 5.3 Training Time and Scalability

Decision Trees and Logistic Regression had the shortest training times, making them suitable for real-time applications or scenarios with limited computational resources. Although Random Forests and Gradient Boosting achieved higher accuracy, their longer training times and memory usage can be prohibitive in resource-constrained environments.

### 5.4 Trade-offs and Hyperparameter Tuning

The study reveals a trade-off between accuracy and computational cost. Gradient Boosting and Random Forests offer high accuracy at the expense of longer training times and memory consumption. Conversely, Decision Trees and Logistic Regression offer quicker, albeit less accurate, solutions. Hyperparameter tuning, such as adjusting the number of trees in Random Forests or the learning rate in Gradient Boosting, significantly impacts model performance. This underscores the importance of hyperparameter optimization in achieving optimal results.

## 6. CONCLUSION

This comprehensive comparative analysis provides a detailed understanding of the performance of various machine learning algorithms on structured datasets. The study's findings suggest that Gradient Boosting and Random Forests are the top-performing algorithms, making them ideal for applications where high accuracy is paramount, such as medical diagnoses or financial risk assessments. On the other hand, Decision Trees and Logistic Regression are more suitable for real-time applications or scenarios with limited computational resources, where speed and efficiency are crucial.

The study highlights the importance of considering multiple factors when selecting an algorithm, including accuracy, computational efficiency, and scalability. By balancing these factors, practitioners can choose the most appropriate algorithm for their

specific use case, ensuring optimal performance and resource utilization.

### 6.1 Future Research Directions

This study provides a foundation for future research in machine learning on structured data. Future

studies can explore additional algorithms, such as neural networks or ensemble methods, to further enhance performance. Advanced techniques, such as transfer learning or feature engineering, can also be investigated to improve the accuracy and efficiency of machine learning models.

## REFERENCES

- [1] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32.
- [2] Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273-297.
- [3] Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232.
- [4] Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Science & Business Media.
- [5] Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1(1), 81-106.
- [6] Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- [7] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785-794). ACM.
- [8] Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- [9] Guyon, I., & Elisseeff, A. (2003). An introduction to variable and feature selection. *Journal of Machine Learning Research*, 3, 1157-1182.
- [10] Ho, T. K. (1995). Random decision forests. In *Proceedings of the 3rd International Conference on Document Analysis and Recognition* (Vol. 1, pp. 278-282). IEEE.
- [11] Dietterich, T. G. (2000). An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine Learning*, 40(2), 139-157.
- [12] Platt, J. C. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Large Margin Classifiers* (pp. 61-74). MIT Press.
- [13] Zhou, Z.-H. (2012). *Ensemble Methods: Foundations and Algorithms*. CRC Press.
- [14] Nguyen, H., & De Baets, B. (2020). Multiclass support vector machines: A literature review. *Computers & Mathematics with Applications*, 77(8), 1989-2004.
- [15] Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence* (Vol. 2, pp. 1137-1143).