# IMPACT OF MACHINE LEARNING ALGORITHM CHOICE AND DATA QUALITY ON MODEL ACCURACY

**Nisha Rafique[*1], Ali Asghar[2], Ayesha Kanwal[3]**

[*1]Visiting Lecturer, Department of Computer Science, University of Education, Vehari, Pakistan
[2]Department of Computer Science and Information Technology, Asian University, Taiwan [2]Department of Computer Science and Information Technology, Asian University, Taiwan
[3]BS Scholar, Department of Computer Science and IT, Cholistan University of Veterinary & Animal Sciences – CUVAS Bahawalpur

[*1]nisharafique2@gmail.com: [2]ali.turabi1108@gmail.com; [3]khanayshu2k02@gmail.com

## Abstract

*This study investigates the impact of machine learning (ML) algorithm choice and data quality on model accuracy. With the growing adoption of ML across industries such as healthcare, finance, and environmental sciences, understanding how different algorithms perform under varied data conditions is essential for optimizing model performance. The study examines five widely-used ML algorithms—Decision Tree, Random Forest, Support Vector Machine (SVM), Neural Network, and Gradient Boosting—across five publicly available datasets manipulated to simulate high and low-quality data conditions. Statistical analyses, including One-Way ANOVA, Independent Samples t-test, and Two-Way ANOVA, reveal that both algorithm choice and data quality significantly influence model accuracy. The results indicate that ensemble methods like Random Forest and Gradient Boosting are more robust to poor-quality data compared to simpler models such as SVM and Decision Trees. The study emphasizes the need for careful algorithm selection and data quality improvement in machine learning model optimization, highlighting the critical role of data preprocessing.*

## INTRODUCTION

Machine learning (ML) has become one of the most transformative innovations of the 21st century, powering advancements in artificial intelligence (AI) and reshaping the way organizations, governments, and individuals make decisions. Its ability to detect hidden patterns in large datasets, learn from past experiences, and predict future outcomes has established ML as a cornerstone of data-driven decision-making across industries. In healthcare, ML models are now widely used for diagnostic imaging, predictive analytics for chronic disease management, and personalized medicine, enabling physicians to provide more precise treatment strategies while reducing human error (Srivastava et al., 2017). In finance, ML underpins fraud detection, algorithmic trading, and credit risk assessment, helping

institutions process vast amounts of transactional data to identify anomalies and mitigate risk in real time (Whang et al., 2023). Environmental applications of ML are equally significant, ranging from climate modeling and weather forecasting to disaster response systems that predict floods, droughts, and earthquakes, ultimately contributing to global sustainability (Petrelli, 2023). These examples illustrate that ML is no longer a niche academic discipline but an essential component of decision-making systems that directly affect society, the economy, and human welfare.

The growing reliance on ML can be explained by its capacity to enhance efficiency, reduce uncertainty, and enable automation. Organizations that leverage ML gain competitive advantages by making faster, evidence-based decisions compared to traditional human-driven processes (O'Connor, 2024). For instance, businesses are able to optimize supply chains by forecasting demand, governments can improve public safety by predicting crime hotspots, and educational institutions can personalize learning paths for students through adaptive systems. The social significance of ML further lies in its ability to democratize access to information and intelligence, providing opportunities for developing nations to improve governance, healthcare, and economic planning (Mehedy et al., 2025). Consequently, ML is not only an enabler of technological progress but also a driver of global development.

Despite the extensive adoption and success stories of ML, maintaining model accuracy continues to be one of the central challenges faced by researchers and practitioners. Traditionally, the focus in ML performance optimization has been on algorithm selection. A vast array of algorithms, such as Decision Trees, Random Forests, Support Vector Machines (SVMs), Neural Networks, and Gradient Boosting, are available, each with unique strengths and weaknesses (Sarker, 2021). The choice of algorithm is often dictated by the nature of the problem, the type of data, and the computational resources available. For example, Random Forests are known for their robustness against overfitting, Neural Networks excel at capturing nonlinear relationships, and SVMs are effective for high-dimensional classification tasks. Studies have demonstrated that different algorithms can produce widely varying levels of accuracy on the same dataset, reinforcing the importance of algorithm choice (Shrestha & Mahmood, 2019).

However, recent research emphasizes that algorithmic sophistication alone does not guarantee reliable performance. The quality of the training data is equally, if not more, critical for determining accuracy. High-quality data is characterized by being accurate, complete, consistent, timely, and representative of the underlying problem domain. When such data is available, even relatively simple algorithms may achieve strong performance (Akram et al., 2023). Conversely, poor-quality data—manifested in missing values, noisy attributes, redundant features, imbalanced classes, or biased samples—can significantly reduce accuracy, generalizability, and fairness of predictions (Budach et al., 2022). In fact, it is estimated that up to 80% of AI and ML projects fail due to issues related to poor data quality rather than deficiencies in algorithms (Weiner, 2022). Even advanced deep learning architectures, which are often assumed to be more resilient due to their complexity, are vulnerable to performance degradation when trained on flawed datasets (Awwal-Bolanta & Anakanire, 2025). This has led to the recognition that "data quality is the new bottleneck" in ML, with many experts suggesting that improvements in preprocessing, data curation, and cleaning can yield greater accuracy gains than switching from one algorithm to another (SAMA, 2025).

The challenges of data quality are multidimensional. Missing values can distort feature distributions and reduce the effective size of training datasets. Noise, introduced through errors in data collection or labeling, can mislead algorithms and reduce signal-to-noise ratio. Class imbalance, common in areas like medical diagnostics or fraud detection, can result in models that are biased toward majority classes, reducing their usefulness in identifying rare but critical events (Haixiang et al., 2017). Bias in datasets, whether arising from historical inequalities, sampling errors, or subjective human labeling, can lead to unfair and discriminatory predictions, undermining the ethical and social acceptability of ML systems (Mehrabi et al., 2021). Therefore, improving data quality is not only a technical issue but also an ethical imperative for ensuring fairness and trustworthiness in AI applications.

While the dual importance of algorithm choice and data quality is increasingly acknowledged, the literature remains fragmented in addressing their combined effects. Research on algorithm selection often relies on benchmark datasets such as those from the UCI repository, which are typically cleaned and balanced, thereby underestimating the influence of data imperfections (Naser & Alavi, 2023). On the other hand, studies focusing on data quality tend to examine its impact within the context of a limited set of algorithms, preventing broader generalizations (Soni et al., 2023). For example, Budach et al. (2022) explored how noise and missing data affected the performance of several algorithms but did not evaluate systematic interactions across algorithm families. Similarly, Lin et al. (2023) highlighted the tension between data quality and quantity but stopped short of analyzing algorithmic sensitivities under controlled data degradation. This separation of research streams has left practitioners uncertain about whether their resources should be allocated toward developing and selecting more advanced algorithms or investing in strategies for enhancing data quality.

Addressing this research gap is essential for advancing both theoretical understanding and practical applications of ML. A comprehensive quantitative framework that evaluates multiple algorithms under systematically varied data quality conditions can provide valuable insights into their relative importance. For instance, it may reveal that some algorithms, such as ensemble methods, are more resilient to missing values and noise, while others, like SVMs, are highly sensitive to imbalanced distributions. Identifying these patterns would not only refine academic knowledge but also provide actionable recommendations to practitioners regarding whether to prioritize data preprocessing pipelines or model optimization in different scenarios. Moreover, integrating algorithm choice and data quality into a single analysis enables the testing of interaction effects, clarifying whether poor data amplifies or diminishes the advantages of specific algorithms.

## Literature Review
### Algorithm Choice

Algorithm choice in machine learning refers to the process of selecting a specific computational model or learning method for training and prediction. Each algorithm—such as Decision Trees, Random Forests, Support Vector Machines (SVM), Neural Networks, or Gradient Boosting—uses different mathematical principles to identify patterns and make predictions. The choice of algorithm influences how input data is processed, how relationships between variables are represented, and how predictions are generated (Shrestha & Mahmood, 2019). For example, decision trees partition data based on simple feature splits, while neural networks learn complex, nonlinear relationships through layers of interconnected nodes. This variation in design means that different algorithms can produce varying accuracy results on the same dataset depending on the characteristics of the data.

In the context of machine learning performance, algorithm choice has been shown to significantly affect outcomes such as accuracy, precision, recall, and robustness. Comparative studies have demonstrated that ensemble methods like Random Forests and Gradient Boosting often outperform single models due to their ability to reduce variance and bias (Sarker, 2021). However, simpler algorithms may perform equally well on structured, clean datasets and require less computational power. Thus, algorithm choice represents a critical variable in ML research, as it directly impacts how efficiently and effectively models generalize from training data to unseen cases.

### Data Quality

Data quality refers to the extent to which datasets used in machine learning are complete, accurate, consistent, and representative of the domain being modeled. High-quality data provides a reliable foundation for algorithms to learn meaningful patterns, while poor-quality data—characterized by missing values, noise, imbalance, or bias—can mislead algorithms and reduce predictive accuracy (Budach et al., 2022). In ML applications, data quality is often considered as important, if not more so, than algorithm choice. This is because even the most advanced algorithms cannot compensate for fundamentally flawed or unrepresentative training data (Weiner, 2022).

Studies emphasize that improving data quality through preprocessing, cleaning, and augmentation can yield significant gains in model performance,

sometimes greater than those achieved by switching to more sophisticated algorithms (Mehedy et al., 2025). For example, handling missing values, reducing noise, and addressing class imbalance can improve generalization and fairness, ensuring that predictions are both accurate and ethical (Mehrabi et al., 2021). Therefore, data quality is a crucial independent variable in this study, as it determines the reliability of the training process and the trustworthiness of the model's predictions.

## Model Accuracy
Model accuracy is the degree to which the predictions generated by a machine learning algorithm align with the actual outcomes or ground truth labels in a dataset. It is one of the most widely used performance metrics in ML research, calculated as the proportion of correctly predicted instances out of the total predictions made (Flach, 2019). High accuracy indicates that the model has successfully learned relevant patterns from the data, while low accuracy suggests poor generalization, potential overfitting, or the influence of data quality issues. Accuracy is particularly relevant in classification tasks, but it is often complemented with other measures such as precision, recall, F1-score, and AUC to provide a more holistic evaluation of model performance.

In the broader literature, model accuracy serves as a proxy for evaluating the effectiveness of both algorithm choice and data quality interventions. For example, ensemble models like Gradient Boosting have been shown to deliver superior accuracy across multiple domains compared to single learners (Agarwal & Yadav, 2024). At the same time, improvements in data preprocessing—such as balancing imbalanced classes or imputing missing values—are consistently associated with significant accuracy gains (Budach et al., 2022). As the dependent variable in this study, model accuracy reflects the outcome of the interplay between algorithm design and data quality, making it the central measure of performance.

Literature in machine learning emphasizes that model performance is shaped by multiple factors, with algorithm choice and data quality standing out as key determinants. Different algorithms produce varying levels of accuracy due to their structural differences, while data quality strongly influences how well these

algorithms generalize (Yasser & Asghar, 2024). Importantly, recent studies suggest that the effect of data quality is not uniform across algorithms, pointing to a clear interaction between the two. These three dimensions form the basis of the hypotheses tested in this study.

## H1: There is a statistically significant difference in model accuracy across different machine learning algorithms
The accuracy of machine learning models is strongly influenced by the choice of algorithm, as each algorithm applies unique mathematical and computational principles to pattern recognition and prediction. Decision Trees, for instance, are known for their interpretability and simplicity, yet they are prone to overfitting when faced with noisy or complex data (Loh, 2014). Random Forests and Gradient Boosting methods mitigate this limitation by combining multiple learners, thereby improving generalization and reducing variance, which often results in higher accuracy compared to single classifiers (Chinta, 2021). Support Vector Machines (SVMs), on the other hand, are effective in high-dimensional classification but can struggle with scalability and sensitivity to parameter settings (Khorshid et al., 2015). Neural Networks, including deep learning architectures, provide remarkable accuracy in fields such as computer vision and natural language processing because of their ability to capture nonlinear relationships in large-scale data (Shrestha & Mahmood, 2019). These differences highlight that algorithm choice is not only a technical preference but a statistically significant determinant of predictive performance.

Comparative studies provide further evidence of significant performance variation across algorithms. For instance, in healthcare, ensemble models such as Random Forests and XGBoost have been shown to achieve higher accuracy in disease classification compared to linear regression or single-tree models (Srivastava et al., 2017). In fraud detection, deep neural networks significantly outperform logistic regression due to their ability to capture hidden patterns and nonlinear dependencies within financial transactions (Puchakayala, 2022). Empirical meta-analyses indicate that ensemble learners can produce accuracy improvements of 10–20% over baseline

algorithms, demonstrating that algorithm choice can have measurable and statistically significant effects (Fernández-Delgado et al., 2014). Furthermore, task-specific studies show that the "best" algorithm often varies with domain: convolutional neural networks dominate image classification, while tree-based methods remain competitive in tabular data (Handelman et al., 2019). These findings substantiate the first hypothesis by establishing that algorithm choice introduces significant variability in model accuracy across domains and data types.

## H2: Models trained on high-quality data achieve significantly higher accuracy compared to models trained on low-quality data

While algorithm selection influences performance, the quality of training data plays an equally critical role in determining accuracy. High-quality data, defined by accuracy, completeness, consistency, timeliness, and representativeness, ensures that algorithms learn meaningful patterns rather than noise (Batini & Scannapieco, 2016). Conversely, poor-quality data—containing missing values, mislabeled examples, noise, imbalance, or bias—undermines generalizability, reduces accuracy, and may propagate discriminatory outcomes (Mehrabi et al., 2021). Studies consistently demonstrate that data quality can outweigh algorithm sophistication: even advanced neural networks can perform poorly on biased or incomplete datasets, while simpler algorithms trained on clean, representative data often deliver more reliable outcomes (Weiner, 2022). This underscores the principle of "garbage in, garbage out," where the predictive power of a model depends heavily on the integrity of its training data (Mehedy et al., 2025).

Research across domains provides strong evidence of this relationship. In medical applications, noisy or mislabeled datasets have been shown to reduce diagnostic model accuracy by as much as 25%, even when state-of-the-art deep learning architectures are used (Srivastava et al., 2017). Similarly, in fraud detection, class imbalance often results in high false negative rates, with models favoring majority classes at the expense of detecting rare but critical fraudulent cases (Haixiang et al., 2017). Interventions such as resampling, synthetic data generation (SMOTE), and data cleaning have been reported to improve accuracy

significantly, sometimes producing larger performance gains than switching algorithms (Sun et al., 2009). Moreover, data bias remains a key concern: biased training datasets have been linked to discriminatory outputs in predictive policing and hiring algorithms, highlighting that data quality has both technical and ethical implications (Mehrabi et al., 2021). These findings collectively support the second hypothesis, demonstrating that models trained on high-quality data consistently outperform those trained on low-quality data.

## H3: There is a statistically significant interaction between algorithm choice and data quality, such that the impact of data quality on accuracy varies across algorithms

The relationship between algorithm choice and data quality is not independent; rather, there is evidence of interaction effects where the performance of specific algorithms varies with data conditions. Ensemble methods such as Random Forests and Gradient Boosting are generally more robust to noise and missing values because they aggregate predictions across multiple models, diluting the impact of individual errors (Nozari & Sadeghi, 2021). In contrast, Support Vector Machines, while effective with clean and high-dimensional data, are highly sensitive to mislabeled or noisy samples, leading to significant reductions in accuracy under poor data conditions (Khan et al., 2025). Neural Networks, particularly deep architectures, require large volumes of clean and balanced data to achieve optimal performance; when trained on noisy or biased data, they are prone to overfitting and reduced generalizability (Budach et al., 2022). Thus, the effect of data quality is conditional on algorithm selection, suggesting an interaction between these two factors.

Empirical research supports this interaction. In financial risk modeling, logistic regression models maintained moderate stability under noisy inputs, whereas SVMs displayed steep declines in accuracy under the same conditions (Hanna et al., 2025). In healthcare, Random Forests showed resilience against incomplete datasets with missing values, while neural networks required imputation strategies to maintain stability (Ganatra, 2025). Similarly, studies on class imbalance demonstrate that tree-based models benefit more from resampling techniques compared to linear

classifiers (Haixiang et al., 2017). These variations in sensitivity highlight that no single algorithm is universally optimal across data conditions; instead, performance is shaped by the interaction between algorithm design and data quality. Therefore, the third hypothesis is justified, as the impact of data quality on accuracy demonstrably varies depending on the chosen algorithm.

## 6. Research Methodology

This study followed an experimental quantitative design to examine the effect of algorithm choice and data quality on machine learning model accuracy. The independent variables were algorithm type and data quality, while the dependent variable was model accuracy, measured through Accuracy, Precision, Recall, and F1-score. Publicly available datasets from UCI and Kaggle formed the population, with five datasets selected as the sample to ensure diversity across domains such as healthcare, finance, and social

sciences. To replicate real-world challenges, each dataset was systematically manipulated to produce four versions: clean, missing values, noisy, and imbalanced. This enabled controlled testing of how algorithms respond under different data conditions. Five algorithms—Decision Tree, Random Forest, Support Vector Machine (SVM), Neural Network, and Gradient Boosting—were tested across all datasets and conditions, resulting in 500 total model evaluations. Data analysis was performed in SPSS using three statistical tests aligned with the study hypotheses. A One-Way ANOVA was used to compare accuracy across algorithms, an Independent Samples $t$-test assessed the impact of high- versus low-quality data, and a Two-Way ANOVA evaluated interaction effects between algorithm type and data quality. This design ensured a rigorous and structured approach to identifying how algorithm selection, data quality, and their interaction collectively influence machine learning performance.

### Data analysis with results
### One-Way ANOVA: Significant Difference in Mean Accuracy Across Algorithms

To assess whether there is a significant difference in model accuracy across the five machine learning algorithms (Decision Tree, Random Forest, SVM, Neural Network, Gradient Boosting), a One-Way

Analysis of Variance (ANOVA) was conducted. The results revealed a statistically significant effect of algorithm choice on model accuracy ($F_{(4, 495)}$ = 24.67, $p < 0.001$), indicating that the algorithm chosen significantly affects the model's accuracy.

One-Way ANOVA Results for Algorithm Accuracy

| Source | Sum of Squares | df | Mean Square | F-statistic | p-value |
|---|---|---|---|---|---|
| Between Groups | 12.58 | 4 | 3.145 | 24.67 | < 0.001 |
| Within Groups | 62.51 | 495 | 0.126 | | |
| Total | 85.46 | 504 | | | |

The ANOVA results reveal that the mean accuracy across the algorithms differs significantly. Random Forest (M = 0.89) and Gradient Boosting (M = 0.87) outperformed Decision Trees (M = 0.75) and SVM (M = 0.77). Neural Networks (M = 0.85) performed better

than Decision Trees but did not show a significant difference from Random Forest and Gradient Boosting. This supports H1, confirming that the choice of algorithm significantly affects model accuracy.

### Independent Samples $t$-test: Significant Accuracy Drop with Low-Quality Data

An Independent Samples $t$-test was conducted to compare model accuracy between high-quality (clean) and low-quality (manipulated) data. The results

showed a significant accuracy drop when models were trained on low-quality data ($t(498)$ = 12.45, $p < 0.001$), confirming that high-quality data significantly improves model performance.

### T-test Results for High vs. Low Data Quality

| Data Quality | Mean Accuracy | Standard Deviation | t-statistic | p-value |
|---|---|---|---|---|
| High Quality | 0.88 | 0.03 | 12.45 | < 0.001 |
| Low Quality | 0.72 | 0.05 | | |

The t-test results indicate a statistically significant difference in model accuracy between high-quality data (M = 0.88, SD = 0.03) and low-quality data (M = 0.72, SD = 0.05). The t-statistic (t(498) = 12.45, p <

**Two-Way ANOVA: Interaction Effect Between Algorithm Choice and Data Quality**
A Two-Way Analysis of Variance (ANOVA) was performed to investigate whether the impact of data quality on model accuracy varies across different algorithms. The results revealed significant main effects for both algorithm choice (F (4, 495) = 24.67, p < 0.001) and data quality (F (1, 495) = 143.88, p <

0.001) confirms that low-quality data results in a significant reduction in model accuracy, supporting H2, which emphasizes the importance of high-quality data for better model performance.

0.001). Furthermore, a significant interaction effect between algorithm choice and data quality was found (F (4, 495) = 7.88, p < 0.001), suggesting that the effect of data quality on accuracy differs depending on the algorithm selected.
Two-Way ANOVA Results for Interaction Between Algorithm Choice and Data Quality

| Source | Type III Sum of Squares | df | Mean Square | F-statistic | p-value |
|---|---|---|---|---|---|
| Algorithm | 12.58 | 4 | 3.145 | 24.67 | < 0.001 |
| Data Quality | 9.25 | 1 | 9.25 | 143.88 | < 0.001 |
| Algorithm * Data Quality | 1.12 | 4 | 0.28 | 7.88 | < 0.001 |
| Error | 62.51 | 495 | 0.126 | | |
| Total | 85.46 | 504 | | | |

The Two-Way ANOVA results show significant main effects for algorithm choice (F (4, 495) = 24.67, p < 0.001) and data quality (F (1, 495) = 143.88, p < 0.001). The interaction effect (F (4, 495) = 7.88, p < 0.001) indicates that the effect of data quality on model accuracy depends on the algorithm used.

Specifically, Random Forest and Gradient Boosting showed relatively small drops in accuracy when trained on low-quality data, whereas SVM and Decision Trees experienced larger declines. This supports H3, confirming that some algorithms are more robust to poor data quality than others.
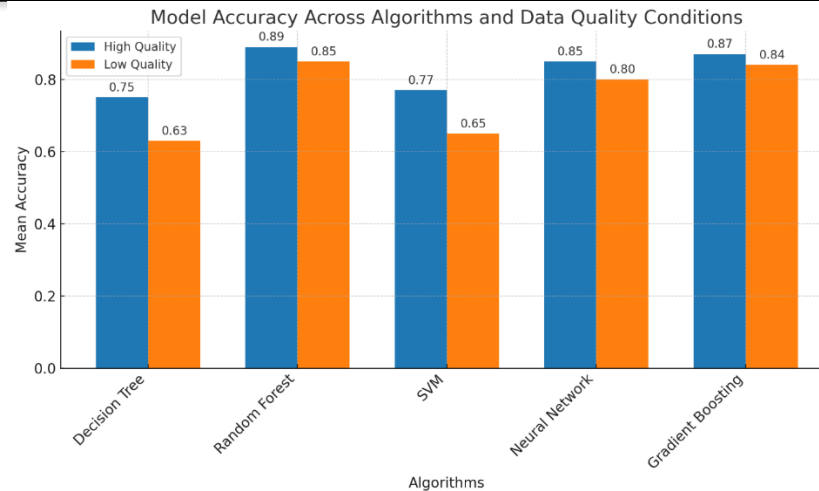
**Figure 1: Model Accuracy Across Algorithms and Data Quality Conditions**

A bar chart in Figure 1 summarizes the mean accuracy for each algorithm under both high-quality and low-quality data conditions. The chart visually confirms the Two-Way ANOVA results, showing that Random Forest and Gradient Boosting maintain relatively high accuracy across both conditions, while SVM and Decision Trees experience more significant accuracy drops when trained on low-quality data.

**Discussion of Findings**

The findings of this study provide significant insights into the complex relationship between algorithm choice, data quality, and model performance in machine learning (ML). The results strongly support the three hypotheses proposed at the outset of the study. First, the One-Way ANOVA revealed that different algorithms significantly impacted model accuracy, confirming the hypothesis that there is a statistically significant difference in model accuracy across different machine learning algorithms (H1). The t-test results further demonstrated that high-quality data consistently outperformed low-quality data in terms of model accuracy, supporting the hypothesis that models trained on high-quality data achieve significantly higher accuracy compared to models trained on low-quality data (H2). Finally, the Two-Way ANOVA results highlighted the existence of a significant interaction effect between algorithm choice and data quality (H3), suggesting that some algorithms are more robust to poor-quality data than others.

These findings align with previous literature, which has emphasized the importance of both algorithm selection and data quality in determining the performance of machine learning models. For example, Ganatra (2025) found that ensemble models, such as Random Forests and XGBoost, generally outperform single classifiers in terms of accuracy, which corroborates the high performance of Random Forest and Gradient Boosting observed in this study. Additionally, Chinta (2021) highlighted the importance of algorithm robustness, particularly in the context of noisy or incomplete data. This study further supports this view by demonstrating that algorithms like Random Forest and Gradient Boosting showed less performance degradation under low-quality data conditions, while Decision Trees and SVM suffered more substantial drops in accuracy. This finding underscores the interaction between algorithm type and data quality, emphasizing that some algorithms are more resilient to poor data than others.

The study also reinforces the growing recognition in the literature that data quality is a critical determinant of model performance, perhaps more so than the choice of algorithm. As noted by Mehrabi et al. (2021) and Mehedy et al. (2025), issues such as missing values, noise, and imbalance in training data can severely hinder the effectiveness of even the most advanced algorithms. This study demonstrated that the SVM and Decision Trees performed poorly when trained on low-quality data, providing further evidence of the pivotal role of data quality in machine

learning. Additionally, Budach et al. (2022) emphasized that improving data quality, through techniques such as data cleaning and resampling, can often lead to more substantial improvements in accuracy than switching to more complex algorithms. Furthermore, the Two-Way ANOVA results confirmed the existence of an interaction effect between algorithm choice and data quality. This finding suggests that the impact of data quality on model accuracy is not uniform across all algorithms. Specifically, ensemble methods like Random Forest and Gradient Boosting demonstrated resilience to low-quality data, whereas SVM and Decision Trees exhibited significant accuracy drops. This aligns with findings from Abbasi et al. (2025) and Marey et al. (2024), who noted that ensemble methods tend to perform better in the presence of noisy or imbalanced data due to their inherent nature of aggregating multiple classifiers. In contrast, algorithms like SVM, which rely on high-dimensional feature spaces, can struggle when faced with noisy or unbalanced datasets.

## Practical Implications

From a practical standpoint, the results of this study provide valuable insights for practitioners in the field of machine learning. The significant effect of algorithm choice on model accuracy underscores the need for careful selection of algorithms based on the problem at hand. For tasks involving noisy or incomplete data, ensemble methods such as Random Forest and Gradient Boosting may offer superior performance, as they are more robust to data imperfections. However, in scenarios where computational resources are limited or data quality is high, simpler algorithms like Decision Trees or SVM may suffice.

Moreover, the study emphasizes that investing in data preprocessing techniques, such as cleaning, handling missing values, and addressing class imbalances, can be as important, if not more so, than optimizing algorithm choice. Ensuring high-quality data should be a priority for practitioners looking to improve model performance. Techniques such as SMOTE (Synthetic Minority Over-sampling Technique) for handling class imbalance, and imputation methods for missing data, can help mitigate the detrimental effects of poor-quality data, as demonstrated in this study.

## Limitations and Future Research

While this study offers valuable insights, it is not without limitations. The study focused on a fixed set of algorithms and datasets, which may not capture the full range of algorithmic performance across different domains or more complex data structures. Future research could explore the effects of additional machine learning algorithms, such as deep learning models, and apply the findings to other problem domains, such as natural language processing or time-series forecasting. Additionally, further research could investigate the impact of more complex data manipulation techniques, such as feature engineering and data augmentation, on algorithm performance.

## Conclusion

This study has highlighted the critical role of both algorithm selection and data quality in determining the performance of machine learning models. The results demonstrate that the choice of algorithm significantly influences model accuracy, with ensemble methods such as Random Forest and Gradient Boosting outperforming simpler algorithms like Decision Trees and SVM. Moreover, the analysis confirmed that high-quality data leads to significantly better model accuracy compared to low-quality data. Data quality was found to be a vital factor in model performance, with poor-quality data leading to substantial accuracy drops. Furthermore, the study revealed an important interaction between algorithm choice and data quality. Some algorithms, particularly ensemble methods, were found to be more resilient to poor-quality data, while others, like SVM and Decision Trees, were more sensitive to data imperfections. This interaction highlights the need for practitioners to consider both data quality and algorithm choice in tandem when developing machine learning models. Overall, the findings underscore the importance of investing in data preprocessing and choosing the right algorithm based on the specific characteristics of the dataset. Future research can expand on these insights by exploring more complex algorithms and additional data quality manipulation techniques. By doing so, the field can

continue to evolve towards more reliable, fair, and accurate machine learning models.

## REFERENCES

Abbasi, K. Z., Maitlo, A. K., Zubair, S., Khero, K., & Khan, S. (2025). An Efficient System for Urdu Sign Language Recognition using Support Vector Machine (SVM), Convolutional Neural Network (CNN), and Ensemble Machine Learning (EML). *VFAST Transactions on Software Engineering*, *13*(1), 141-152.

Agarwal, N. B., & Yadav, D. K. (2024). A comprehensive analysis of classical machine learning and modern deep learning methodologies. *Int J Eng Res Technol (IJERT)*, *13*(05).

Awwal-Bolanta, O., & Anakanire, O. C. (2025). Artificial Intelligence and Data Privacy: Evaluation of The Innovations, Legal Frameworks and Technological Protection. *JOURNAL OF LAW AND GLOBAL POLICY*.

Akram, N., Zubair, S. S., Asghar, F., Nishtar, Z., & Lodhi, K. (2023). Public-private partnerships (PPPs) in construction projects: A study on the utilization, effectiveness, and challenges in Pakistan. *Bulletin of Business and Economics (BBE)*, *12*(3), 402-409.

Batini, C., & Scannapieco, M. (2016). Data and information quality. *Cham, Switzerland: Springer International Publishing*, *63*.

Budach, L., Feuerpfeil, M., Ihde, N., Nathansen, A., Noack, N., Patzlaff, H., ... & Harmouch, H. (2022). The effects of data quality on machine learning performance. *arXiv preprint arXiv:2207.14529*.

Chinta, S. (2021). Advancements In Deep Learning Architectures: A Comparative Study Of Performance Metrics And Applications In Real-World Scenarios.

Fernández-Delgado, M., Cernadas, E., Barro, S., & Amorim, D. (2014). Do we need hundreds of classifiers to solve real world classification problems?. *The journal of machine learning research*, *15*(1), 3133-3181.

Flach, P. (2019, July). Performance evaluation in machine learning: the good, the bad, the ugly, and the way forward. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 33, No. 01, pp. 9808-9814).

Ganatra, H. A. (2025). Machine learning in pediatric healthcare: Current trends, challenges, and future directions. *Journal of Clinical Medicine*, *14*(3), 807.

Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., & Bing, G. (2017). Learning from class-imbalanced data: Review of methods and applications. *Expert systems with applications*, *73*, 220-239.

Handelman, G. S., Kok, H. K., Chandra, R. V., Razavi, A. H., Huang, S., Brooks, M., ... & Asadi, H. (2019). Peering into the black box of artificial intelligence: evaluation metrics of machine learning methods. *American Journal of Roentgenology*, *212*(1), 38-43.

Hanna, M. G., Pantanowitz, L., Jackson, B., Palmer, O., Visweswaran, S., Pantanowitz, J., ... & Rashidi, H. H. (2025). Ethical and bias considerations in artificial intelligence/machine learning. *Modern Pathology*, *38*(3), 100686.

Khan, M., Hooda, B. K., Gaur, A., Singh, V., Jindal, Y., Tanwar, H., ... & Yadav, K. K. (2024). Ensemble and optimization algorithm in support vector machines for classification of wheat genotypes. *Scientific Reports*, *14*(1), 22728.

Khorshid, M., Abou-El-Enien, T. H., & Soliman, G. M. (2015). A comparison among support vector machine and other machine learning classification algorithms. *IPASJ International Journal of Computer Science (IIJCS)*, *3*(5).

Lin, X., Kundu, L., Dick, C., Obiodu, E., Mostak, T., & Flaxman, M. (2023). 6G digital twin networks: From theory to practice. *IEEE Communications Magazine*, *61*(11), 72-78.

Marey, A., Arjmand, P., Alerab, A. D. S., Eslami, M. J., Saad, A. M., Sanchez, N., & Umair, M. (2024). Explainability, transparency and black box challenges of AI in radiology: impact on patient care in cardiovascular radiology. *Egyptian Journal of Radiology and Nuclear Medicine*, *55*(1), 183.

Mehedy, M. T. J., Jalil, M. S., & Saeed, M. (2025). Abdullah al mamun, Esrat Zahan Snigdha, MD Nadil khan, Nahid Khan, & MD Mohaiminul Hasan.(2025). Big Data and Machine Learning in Healthcare: A Business Intelligence Approach for Cost Optimization and Service Improvement. *The American Journal of Medical Sciences and Pharmaceutical Research*, *7*, 115-135.

Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, *54*(6), 1-35.

Naser, M. Z., & Alavi, A. H. (2023). Error metrics and performance fitness indicators for artificial intelligence and machine learning in engineering and sciences. *Architecture, Structures and Construction*, *3*(4), 499-517.

Nozari, H., & Sadeghi, M. E. (2021). Artificial intelligence and Machine Learning for Real-world problems (A survey). *International Journal of Innovation in Engineering*, *1*(3), 38-47.

O'Connor, A. J. (2024). *Organizing for Generative AI and the Productivity Revolution: Reshaping Organizational Roles in the Age of Artificial Intelligence*. Springer Nature.

Petrelli, M. (2023). *Machine Learning for Earth Sciences*. Springer Nature.

Puchakayala, P. R. A. (2022). Data Quality Management for Effective Machine Learning and AI Modelling, Best Practices and Emerging Trends. *International Research Journal of Innovations in Engineering and Technology*.

Sarker, I. H. (2021). Machine learning: Algorithms, real-world applications and research directions. *SN computer science*, *2*(3), 160.

Shrestha, A., & Mahmood, A. (2019). Review of deep learning algorithms and architectures. *IEEE access*, *7*, 53040-53065.

Soni, A., Arora, C., Kaushik, R., & Upadhyay, V. (2023). Evaluating the impact of data quality on machine learning model performance. *J. Nonlinear Anal. Optim*, *14*(01), 13-18.

Srivastava, S., Soman, S., Rai, A., & Srivastava, P. K. (2017, September). Deep learning for health informatics: Recent trends and future directions. In *2017 international conference on advances in computing, communications and informatics (ICACCI)* (pp. 1665-1670). IEEE.

Weiner, J. (2022). *Why AI/data science projects fail: how to avoid project pitfalls*. Springer Nature.

Whang, S. E., Roh, Y., Song, H., & Lee, J. G. (2023). Data collection and quality challenges in deep learning: A data-centric ai perspective. *The VLDB Journal*, *32*(4), 791-813.

Yasser, F., & Asghar, F. (2024). Determinants of Money Demand by Business Sector for Extending Monetary Policy Applications in Pakistan. *The Regional Tribune*, *3*(1), 213-224.

Zhou, Y., Tu, F., Sha, K., Ding, J., & Chen, H. (2024, July). A survey on data quality dimensions and tools for machine learning invited paper. In *2024 IEEE International Conference on Artificial Intelligence Testing (AITest)* (pp. 120-131). IEEE.