# Comparative Analysis of Machine Learning Models for Classifying Human DNA Sequences: Performance Metrics and Strategic Recommendations

**Gregorius Airlangga**

Engineering Faculty, Information System Study Program, Atma Jaya Catholic University of Indonesia, Jakarta, Indonesia
Correspondence Author Email: gregorius.airlangga@atmajaya.ac.id

**Abstract**−This study presents a comprehensive evaluation of seven machine learning models applied to the classification of human DNA sequences, highlighting their performance and potential applications in genomics. We explored Logistic Regression, Support Vector Machines (SVM), Random Forest, Decision Trees, Gradient Boosting, Naive Bayes, and XGBoost, using a 5-fold StratifiedKFold cross-validation method to ensure robustness and reliability in our findings. Naive Bayes demonstrated exceptional performance with near-perfect accuracy, precision, recall, and F1 scores, suggesting its suitability for rapid and efficient genomic classification. Logistic Regression also showed high efficacy, proving effective even in multi-class classifications of complex genetic data. Conversely, Decision Trees and SVM struggled with overfitting and computational efficiency, respectively, indicating the need for careful parameter tuning and optimization in practical applications. The study addresses these challenges and proposes strategies for enhancing model robustness and computational efficiency, such as advanced regularization techniques and hybrid modeling approaches. These insights not only aid in selecting appropriate models for specific genomic tasks but also pave the way for future research into integrating machine learning with genomic science to advance personalized medicine and genetic research. The findings encourage ongoing refinement of these models to unlock further potential in genomic applications.

**Keywords**: Machine Learning; Genomic Classification; DNA Sequencing; Comparative Analysis; Model Optimization

## 1. INTRODUCTION

The advent of high-throughput genomic sequencing technologies has heralded a new era in biological sciences, significantly influencing fields such as genomics, molecular biology, and medical diagnostics [1]–[3]. These technologies enable the generation of vast datasets that provide unprecedented insights into the genetic bases of diseases, individual variability in drug response, and the intricate mechanisms of life at a molecular level [4]–[6]. However, the benefit of this data can only be realized through effective analytical tools capable of parsing and making sense of the sequences, which are often vast and complex [7]. The use of machine learning (ML) in bioinformatics is not a new concept but has seen exponential growth in relevance and application due to the rise of genomic big data [8]. Initially, bioinformatics relied heavily on alignment-based methods and simple statistical techniques for tasks such as sequence classification and motif finding [9]. However, as genome sequencing has become cheaper and more commonplace, the data generated has grown exponentially in size and complexity, surpassing the capabilities of traditional methods. Recent literature, including comprehensive reviews by [10] and empirical studies by [11], underscores the shift towards machine learning models. These models, particularly ensemble methods and advanced algorithms like deep learning, have demonstrated a significant increase in accuracy and efficiency in detecting complex patterns within large-scale genomic datasets [12].

The rapid evolution of sequencing technologies has not only made genomic data abundant but also increasingly complex. The urgency to develop advanced computational approaches is underscored by the need to understand this complexity in a manner that is scientifically reliable and clinically actionable [13]–[15]. Diseases such as cancer are known to have a genetic component that, when understood, can lead to more effective personalized treatments [16]. Hence, there is a pressing need for robust ML models that can provide insights into genetic variations and their consequences, facilitating advances in personalized medicine [17]. Current research at the intersection of ML and genomics involves a diverse array of methodologies, from traditional algorithms like Random Forests and SVMs to more sophisticated models such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) [18]. Each of these models offers different strengths; for instance, CNNs are particularly adept at pattern recognition in sequence data, making them suitable for tasks like identifying regulatory motifs [19]. Conversely, SVMs and Random Forests are valued for their interpretability and robustness in smaller, less complex datasets. This research focuses on a balanced evaluation of both traditional and contemporary models, assessing them across multiple performance metrics to identify the most effective approaches for specific genomic classification tasks [20].

This study aims to systematically evaluate a range of machine learning models to classify human DNA sequences into one of several predefined classes. By doing so, it seeks to identify which models provide the best performance in terms of accuracy, efficiency, and interpretability. The overarching goal is to recommend models that could potentially improve the speed and accuracy of genetic disease diagnosis and prognosis in clinical settings. Although there is a wealth of research on applying ML models to genomic data, there remains a significant deficiency in studies that perform head-to-head comparisons of a broad spectrum of models on the

same dataset. Many studies focus on a single model type or a small group of similar models, which does not provide the comparative insight needed to guide model selection in practical applications. Moreover, there is a notable lack of comprehensive analysis regarding the trade-offs between model complexity and interpretability, particularly in how these factors impact the deployment of models in different genomic research scenarios.

This paper contributes to the field by providing a detailed comparative analysis of seven distinct machine learning models applied to the classification of human DNA sequences. It offers a nuanced discussion on the strengths and weaknesses of each model, providing guidance for their application in both research and clinical settings. Additionally, this study introduces a methodological framework for evaluating and comparing the efficacy of different ML models on genomic data, which could be adopted in future research. The remainder of this article is organized as follows: Section 2 outlines the methodology, detailing the data preparation, model implementation, experimental setup, and evaluation criteria. Section 3 presents a comprehensive analysis of the results, providing a detailed comparison of model performances. Section 4 discusses these results in the context of current bioinformatics challenges and the potential implications for future genomic research. Finally, Section 5 concludes with a summary of the findings, highlighting key takeaways and suggesting directions for future research.

# 2. RESEARCH METHODOLOGY

This study adopts a comprehensive approach to evaluate and compare the efficacy of various machine learning (ML) models in classifying human DNA sequences into predefined classes. The methodology is structured to ensure reproducibility and robust evaluation, encompassing data preparation, model implementation, cross-validation procedures, and detailed performance metrics as presented in the figure 1.
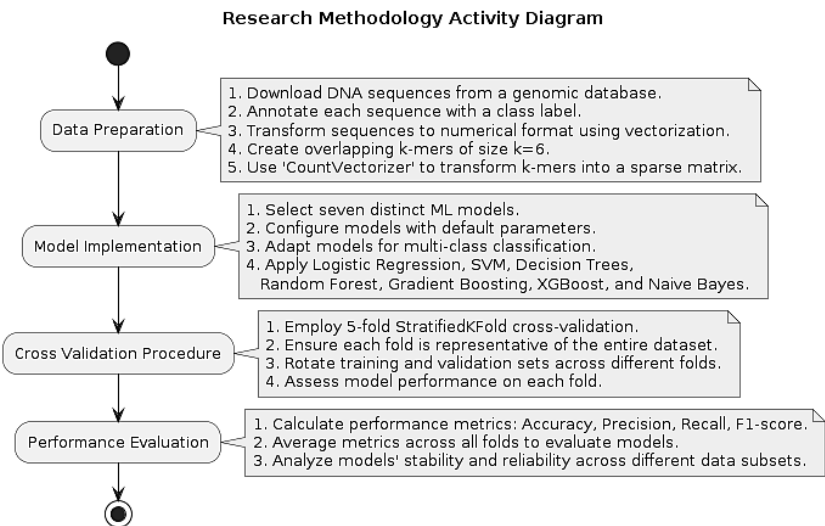


**Figure 1.** Research Methodology Activity Diagram

## 2.1 Data Preparation

The dataset central to our research comprises human DNA sequences that have been publicly sourced from a genomic database and can be downloaded from [21]. Each sequence within the dataset is annotated with a class label, indicating the sequence's association with particular biological functions or characteristics. Overall, the dataset encompasses seven distinct classes. The preliminary phase of data handling involves the transformation of these sequences into a format amenable to ML modeling. This transformation, commonly referred to as vectorization, converts the complex, raw DNA sequences into a numerical format that can be effectively processed by various ML algorithms. Vectorization in this context is achieved through the creation of overlapping k-mers from the DNA sequences. Specifically, each sequence is broken down into contiguous subsequences, each of which is k nucleotides long, where k is set to 6. This value was selected based on initial exploratory analyses which determined that k=6 offers an optimal balance between computational efficiency and the ability of the models to perform predictive tasks effectively. These k-mers serve as individual features within the ML models.

Once the sequences are decomposed into k-mers, we employ a tool known as 'CountVectorizer'. This tool is pivotal in transforming the k-mers into a high-dimensional sparse matrix. In this matrix, each column uniquely corresponds to one k-mer, and each row represents a DNA sequence from our dataset. The values within the matrix cells denote the frequency of each k-mer within a given sequence, providing a structured way to quantify the presence and prevalence of specific nucleotide patterns across the dataset. This method of data preparation

not only preserves the essential biological information inherent in the sequences but also transforms it into a statistically analyzable format. By doing so, we facilitate the application of advanced statistical models and machine learning techniques to uncover patterns and associations that may be indicative of underlying biological processes or phenotypic characteristics associated with the DNA sequences. The transition from raw data to a structured, numerical format is critical in enabling the subsequent phases of model implementation and validation, ensuring that the analyses conducted are both meaningful and scientifically rigorous.

## 2.2 Model Implementation

In this study, the implementation of machine learning models is a core component, aimed at evaluating the capabilities of various algorithms in classifying human DNA sequences. Seven distinct models were carefully selected, each recognized for its applicative strengths and prevalence in genomic studies. This diverse selection enables a thorough exploration of the tools available for genomic classification, facilitating a comprehensive understanding of how different algorithms perform under the same experimental conditions. Logistic Regression and Support Vector Machine (SVM) are foundational tools in the realm of machine learning, known primarily for their robust performance in binary classification tasks. For the purposes of this study, both models were adeptly adapted to handle multi-class classification challenges. This adaptation is crucial in genomic applications where the classification task often involves distinguishing between multiple types of genetic sequences, each potentially linked to different biological functions or diseases.

Random Forest and Decision Trees are included in the analysis due to their methodological transparency, which provides clear, intuitive decision rules. This characteristic is particularly valuable in genomic studies, where understanding which features—here, specific k-mers—are most influential in the classification process can provide insights into the biological significance of these genomic segments. These models not only classify but also help elucidate the feature importance, thereby contributing to a deeper biological interpretation. Furthermore, the study incorporates ensemble methods such as Gradient Boosting and XGBoost. These techniques enhance prediction accuracy by aggregating the outputs of multiple weaker models to form a more robust prediction model. Such methods are beneficial in handling the complex and often noisy data characteristic of genomic datasets, where the integration of multiple decision paths can lead to a significant increase in predictive performance. Naive Bayes was selected for its efficiency in handling high-dimensional data, drawing parallels to scenarios commonly encountered in text classification. Its inclusion is strategic, given the high-dimensional nature of the transformed DNA sequence data, where each k-mer represents a dimension in the feature space. Naive Bayes, with its probabilistic approach, is adept at managing the complexities arising from the vast feature spaces, making rapid classifications even when faced with extensive data.

For the practical implementation of these models, the Python library Scikit-Learn was predominantly used due to its extensive suite of machine learning tools and its user-friendly interface, which is particularly suited for academic and research settings. The exception was XGBoost, which necessitated the use of its dedicated library to leverage specific optimizations and functionalities not available in Scikit-Learn. Initially, all models were configured with their default parameters to establish a baseline of performance. Subsequent tuning and adjustments were made based on preliminary results, allowing for optimization tailored to the specific characteristics of the genomic data being analyzed. This careful and detailed approach to model implementation ensures that the study not only assesses the raw predictive power of each algorithm but also considers their practical applicability and efficiency in real-world genomic classification scenarios. The models are evaluated not just in isolation but as part of a broader system of computational tools available for tackling complex biological data.

### 2.2.1 Logistic Regression

For multi-class classification challenges, logistic regression can be extended beyond the binary setting using the softmax function, also known as multinomial logistic regression. The probability that an instance $(x_i)$ belongs to class $(k)$ is modeled by the softmax function $[P(y = k \mid x_i) = \frac{\exp(x_i^T \beta_k)}{\sum_{j=1}^{K} \exp(x_i^T \beta_j)},]$ where $(\beta_k)$ represents the parameter vector associated with class $(k)$, and $(K)$ is the total number of classes. This function ensures that the probabilities sum to 1 over all classes for any given input $(x_i)$, providing a probabilistic framework for multi-class classification. The parameters $(\beta_k)$ are typically learned by maximizing the likelihood of the training data, which equivalently minimizes the cross-entropy between the predicted and actual distributions. This optimization problem can be solved using gradient-based methods, providing a robust framework for handling multi-class classification tasks.

### 2.2.2 Support Vector Machine (SVM)

SVM can be adapted to multi-class classification using strategies such as one-vs-all (OvA) and one-vs-one (OvO). These strategies allow SVM, originally designed for binary classification, to handle multiple classes. In the OvA strategy, a separate binary classifier is trained for each class to distinguish instances of that class from instances of all other classes. The decision function for a class $(k)$ is given by $f_k(x_i) = x_i^T \beta_k + b_k$. The class

with the highest decision value is then selected as the prediction $\hat{y}_i = \arg\max_k f_k(x_i)$. Meanwhile, in the OvO strategy, a binary classifier is trained for every pair of classes. If there are $(K)$ classes, this results in $(\frac{K(K-1)}{2})$ classifiers. Each classifier decides between two classes, and the final class prediction is typically made by a voting scheme among all classifiers. Both of these strategies extend the powerful margin-based framework of SVM to multi-class problems, allowing for clear, robust decision boundaries between classes. SVM is particularly noted for its effectiveness in high-dimensional spaces, making it well-suited for genomic classification tasks where features might represent complex genetic patterns.

### 2.2.3 Decision Tree

Decision Trees are a non-parametric supervised learning method used for classification and regression tasks. The model predicts the value of a target variable by learning simple decision rules inferred from the data features. A decision tree divides the feature space into a number of regions $(R_m)$, and for each region, a simple prediction model is used, typically the mode (for classification) or mean (for regression) of the target variable in that region. Mathematically, a decision tree function $(f)$ for classification can be represented as $f(x_i) = c_m$ if $x_i \in R_m$, where $(x_i)$ is an input feature vector, $(R_m)$ is a region in the feature space determined by the tree's splitting rules, and $(c_m)$ is the class prediction for region $(R_m)$. The regions $(R_m)$ are formed by recursively partitioning the feature space, splitting it at values of one or more features to maximize the homogeneity of the target variable within each resulting region. There are several steps in order to create construction of a Decision Tree, firstly, starting at the tree root, the feature and threshold that yield the largest information gain are chosen to split the data into two child nodes. Then, this process is repeated recursively for each child node. After that, the recursion is completed when a stopping criterion is met (e.g., maximum depth, minimum samples per leaf, or no further improvement). Therefore, at each leaf node, the most common class of training examples sorted into that node is chosen as the prediction $(c_m)$.

### 2.2.4 Random Forest

Random Forest is an ensemble learning method for classification, regression, and other tasks that operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random Forest improves upon the variance of single decision trees by averaging multiple trees that individually suffer from high variance and may overfit the data. The ensemble approach of Random Forests mitigates these issues by building each tree on a different bootstrap sample of the data and using a random subset of features for splitting at each node, which increases diversity in the ensemble and leads to more robust overall predictions. Mathematically, a Random Forest model can be described as $f(x_i) = \frac{1}{M}\sum_{m=1}^{M} f_m(x_i)$, where $(M)$ is the number of trees, $(f_m)$ is the prediction function of the $(m)$-th tree, and the final prediction is obtained by averaging the predictions of all trees for regression or by majority voting for classification. There are benefits of random forest, firstly, robustness, by averaging multiple trees, the model is less sensitive to outliers and noise in the dataset. Secondly, performance, often provides high predictive performance that can compete with the best supervised learning algorithms. Lastly, feature importance, it provides insights into which features are more important in predicting the target variable.

### 2.2.5 Gradient Boosting

Gradient Boosting is a powerful machine learning technique that builds models incrementally in the form of an ensemble of weak prediction models, typically decision trees. The core idea is to construct the new predictor to be maximally correlated with the negative gradient of the loss function, associated with the whole ensemble. The method involves sequentially adding to an ensemble, where each new model incrementally reduces the loss function (usually a differentiable loss function) used to measure the difference between the predicted and actual values. The ensemble prediction $(f)$ is represented as $f(x_i) = \sum_{m=1}^{M} \gamma_m h_m(x_i)$, where $(M)$ is the number of models in the ensemble, $(h_m)$ represents the weak learner added in the $(m)-th$ iteration, and $(\gamma_m)$ is the weight of that learner, calculated to minimize the loss when combined with the previous learners. Gradient Boosting is a machine learning technique that builds models incrementally, forming an ensemble of weak prediction models, typically decision trees. The method starts by initializing the model with a constant predictor, $(f_0(x))$. Then, for each iteration $(m)$ from 1 to $(M)$, it performs the following steps: compute the residual errors from the previous model, fit a new model $(h_m)$ to these residuals, find the coefficient $(\gamma_m)$ that minimizes the loss when $(h_m)$ is added to the ensemble, and update the ensemble model by adding the new model weighted by $(\gamma_m)$. The final model $(f_M(x))$ represents the accumulated contributions of all iterations.

### 2.2.6 XGBoost

XGBoost (Extreme Gradient Boosting) is an optimized gradient boosting library that provides a parallel tree boosting solution. It enhances traditional gradient boosting through regularization terms in the objective function to prevent overfitting and improves handling of sparse data. XGBoost supports various objective functions,

including regression, classification, and ranking. The ensemble model in XGBoost is formed by summing the contributions of ( $M$ ) models $f(x_i) = \sum_{m=1}^{M} \gamma_m h_m(x_i)$, where ( $h_m$ ) are the weak learners and ($\gamma_m$) are their respective weights. Unique features of XGBoost include column sampling by tree, tree pruning, and robust handling of missing values, which enhance the model's scalability and flexibility. Users can also define custom optimization objectives and evaluation criteria, making XGBoost suitable for a wide range of data science problems.

### 2.2.7 Naive Bayes

The Naive Bayes classifier applies Bayes' theorem with the assumption that all features are independent from each other. This assumption simplifies the computation of the conditional probabilities as $P(y = k \mid x_i) = \frac{P(x_i|y=k)P(y=k)}{P(x_i)}$, where ($P(x_i \mid y = k)$) is assumed to follow a specific distribution, such as Gaussian for continuous data or Multinomial for discrete data. Despite its simplicity and the strong independence assumption, Naive Bayes can perform remarkably well in many complex real-world situations.

The likelihood ($P(x_i \mid y = k)$) is modeled using probability distributions appropriate to the type of data in the features. For continuous features, the Gaussian distribution is typically assumed, where the likelihood of each feature is estimated as the probability density of the normal distribution. In cases involving categorical or count data, such as in text classification, the Multinomial distribution is used, where features may represent the frequency of words. Despite the simplicity and stringent independence assumption, Naive Bayes classifiers often perform surprisingly well, particularly in high-dimensional settings like text classification or genomic data analysis. The model's efficiency in managing large datasets with many features is a significant advantage, offering robust performance across various applications. Naive Bayes classifiers are valued for their speed and ease of implementation, providing reliable outcomes even with the large-scale and complex datasets often encountered in real-world scenarios. They are notably effective at managing missing data and require only a minimal amount of training data to estimate the necessary probabilities, making them an appealing choice for many practical situations.

### 2.3 Cross Validation Procedure

In this study, In the exploration of machine learning models for genomic classification, the validation of these models is as crucial as their implementation. To ensure a rigorous and fair comparison of the models' performance, this study employed a 5-fold StratifiedKFold cross-validation technique. This method of validation is particularly suited to scenarios where the data might not be uniformly distributed across different classes, which is often the case in genomic datasets where some genetic sequences are more prevalent than others. The choice of StratifiedKFold cross-validation addresses several critical needs in the evaluation process. Primarily, it ensures that each fold of the data is representative of the entire dataset. StratifiedKFold cross-validation ensures each fold of data is representative of the entire dataset. This stratification maintains the percentage of samples for each class consistent with the overall dataset. Mathematically, this can be formulated by $P(y_i = k \mid \text{Fold } f) = P(y_i = k) \quad \forall k \in \{1,2,...,K\}$, where ($y_i$) represents the class label of the ($i$)-th sample, ($k$) denotes a class, ($K$) is the total number of classes, and ($f$) indicates a specific fold. This is achieved by distributing samples so that each fold contains approximately the same percentage of samples of each target class as the original dataset. Such stratification is essential in avoiding bias in the training process and in the evaluation of the model's performance, as it prevents any single class from dominating the learning process, which could lead to skewed results that do not generalize well across the spectrum of possible inputs. The dataset is split into 5 distinct subsets (folds). Each fold acts once as a validation set while the remaining four folds form the training set. This cycle is repeated 5 times, with each fold serving as the validation set exactly once. The validation process for each fold is defined by Validation Set $V_f = D_f$, Training Set $T_f = D \setminus D_f$.

Stratification plays a pivotal role in maintaining the integrity and validity of the machine learning process, particularly in the field of genomics. In genomic studies, the consequences of class imbalance can be profound, as models trained on imbalanced data may perform well on majority classes but poorly on minority classes, which are often of significant scientific and clinical interest. By ensuring a balanced representation of classes in each fold, StratifiedKFold cross-validation helps in creating more robust models that are less likely to exhibit bias towards the more frequent classes. Moreover, the use of cross-validation in this manner aids in mitigating the risk of overfitting. Overfitting occurs when a model is excessively complex, capturing noise in the training data as if it were true signal, which impairs its performance on new, unseen data. By rotating the training and validation sets across different folds of the data, and ensuring that each instance of the dataset is used for both training and validation in different iterations, the process provides a comprehensive assessment of how well a model is likely to perform in practical scenarios beyond the confines of the experimental setup. The robustness of the models is further evaluated through this cross-validation procedure, which provides not only an estimate of the overall effectiveness of each model but also insights into their stability and reliability across different subsets of the data. This thorough testing framework is crucial for drawing reliable conclusions about the models' capabilities and for ensuring that the models can generalize well to new data, an essential quality for models intended to be used in ongoing genomic research and clinical applications.

## 2.4 Performance Evaluation

In assessing the performance of machine learning models for genomic classification, it is imperative to apply a set of metrics that provide a multifaceted view of each model's capabilities. This study utilizes a comprehensive set of performance metrics: accuracy, precision, recall, and the F1 score, each of which contributes a unique perspective on the effectiveness of the models. Accuracy is the most straightforward of these metrics, offering a high-level view of the model's overall effectiveness across all classes. It quantifies the proportion of total predictions that were correctly classified, giving a general sense of how often the model is correct in its classifications. However, while accuracy is useful for providing a snapshot of performance, it can sometimes be misleading, especially in datasets where some classes are much more prevalent than others. To address the nuances and potential imbalances in the dataset, precision and recall are employed as more discerning metrics. Precision measures the accuracy of positive predictions for each class, which is the ratio of true positive predictions to the total number of instances predicted as positive. This metric is crucial when the consequences of false positives are significant, such as in medical diagnostics where a false positive might lead to unnecessary treatment.

Recall, on the other hand, measures the model's ability to correctly identify all positive samples from the actual positives available within the data. This metric is particularly important in medical or biological contexts where missing a positive instance (a false negative) can be more detrimental than a false positive. For example, failing to identify a genetic marker associated with a disease could prevent a patient from receiving necessary medical attention. The F1 score, which is the harmonic mean of precision and recall, is used to balance these two metrics. It is particularly useful when the cost of false positives and false negatives are equally concerning, providing a single metric that balances both aspects of the model's performance. The F1 score is especially relevant in genomic studies, where both types of errors can have significant implications. For each fold in the 5-fold StratifiedKFold cross-validation, these metrics were calculated to evaluate the models' performance. By averaging these metrics across all folds, the study provides a robust measure of each model's average performance as well as insights into the consistency of this performance across different data subsets. This multi-metric evaluation framework is critical in genomic settings, allowing researchers to understand not only the general effectiveness of each model but also how its performance might vary in different practical scenarios. This approach ensures a thorough understanding of each model's strengths and weaknesses, facilitating informed decisions about their application in real-world genomic classification tasks. The equations for calculating accuracy, precision, recall and F1 score is presented in the equation (1) – (4).

$$Accuracy = \frac{\text{number of correct predictions}}{\text{total predictions}} \tag{1}$$

$$\text{Precision} = \frac{TP}{TP+FP} \tag{2}$$

$$\text{Recall} = \frac{TP}{TP+FN} \tag{3}$$

$$\text{F1 Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \tag{4}$$

Where, $TP$ is true positives, $FP$ is false positives, $FN$ is false negatives. In addition, $TPR(t)$ is the true positive rate at threshold $t$, and $FPR(t)$ is the false positive rate.

# 3. RESULT AND DISCUSSION

## 3.1 Result

This section presents and discusses the findings derived from evaluating seven different machine learning models on human DNA sequence classification using a 5-fold StratifiedKFold cross-validation approach. The primary indicators of each model's effectiveness include accuracy, precision, recall, and F1 score, providing a comprehensive view of their performance in genomic applications. As presented in the table 1, Naive Bayes emerged as the standout performer, demonstrating the highest efficiency across all evaluated metrics with an accuracy of 97.47%, precision of 97.56%, recall of 97.47%, and an F1 score of 97.48%. These results highlight Naive Bayes as a highly effective model for handling the complex, high-dimensional data typical of genomic sequences. The model benefits from its assumption of feature independence, which simplifies computational processes and enhances its suitability for genomic data. Logistic Regression also displayed robust performance, securing an accuracy of 93.63%, precision of 94.50%, recall of 93.63%, and an F1 score of 93.70%. This model proves that even simpler techniques can provide substantial accuracy in multi-class classification tasks, reinforcing its utility in genomic studies where interpretability and performance are both crucial.

XGBoost achieved commendable results with its ensemble approach, yielding an accuracy of 89.13%, precision of 90.45%, recall of 89.13%, and an F1 score of 89.20%. This model's ability to integrate multiple decision trees to form a strong predictive model makes it suitable for datasets where complex model

relationships are a factor. Random Forest performed well with an accuracy of 90.82%, precision of 92.00%, recall of 90.82%, and an F1 score of 90.95%. It excels in handling feature interactions and provides valuable insights into which features are most influential in classification, a useful trait for exploratory genomic analysis. Gradient Boosting, with an accuracy of 83.77%, precision of 87.63%, recall of 83.77%, and an F1 score of 84.12%, showed it could manage the non-linear relationships within genomic data, although not as effectively as some other models. Support Vector Machine (SVM) and Decision Trees reported lower performance metrics compared to the other models. SVM, with an accuracy of 81.94%, precision of 88.39%, recall of 81.94%, and an F1 score of 82.53%, and Decision Trees, with an accuracy of 81.28%, precision of 83.29%, recall of 81.28%, and an F1 score of 81.76%, faced challenges likely due to their sensitivity to parameter settings and susceptibility to overfitting, respectively.

**Table 1.** The model's performance results

| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| Logistic Regression | 0.9363 | 0.9450 | 0.9363 | 0.9370 |
| SVM | 0.8194 | 0.8839 | 0.8194 | 0.8253 |
| Random Forest | 0.9082 | 0.9200 | 0.9082 | 0.9095 |
| Decision Trees | 0.8128 | 0.8329 | 0.8128 | 0.8176 |
| Gradient Boosting | 0.8377 | 0.8763 | 0.8377 | 0.8412 |
| Naive Bayes | 0.9747 | 0.9756 | 0.9747 | 0.9748 |
| XGBoost | 0.8913 | 0.9045 | 0.8913 | 0.8920 |

These results illustrate a significant variability in performance, which can be attributed to each model's inherent characteristics and their interaction with the genomic data structure. Naive Bayes is particularly notable for scenarios where rapid classification is needed. However, its assumption of feature independence may not hold in all genomic datasets, which suggests that its effectiveness could vary. Both Logistic Regression and Random Forest offer a good balance between accuracy and the ability to interpret the model outputs, crucial for understanding the biological implications of genomic classifications. The study highlights the importance of model tuning and feature selection, especially in improving models like SVM, and stresses the need for careful selection based on specific research needs, whether prioritizing accuracy or minimizing false negatives is more critical. This comparative analysis serves as a foundation for selecting suitable machine learning models for genomic classification, advocating for a balanced approach where performance metrics and biological context are both considered. Future research could explore combining the strengths of these models through hybrid or advanced ensemble techniques to further enhance predictive accuracy and reliability in genomic classification.

**3.2 Trade-off Analysis**

In genomic studies, the selection of an appropriate machine learning model involves considering various trade-offs between accuracy, interpretability, computational efficiency, and robustness. The findings from this study highlight the complex interactions between these factors and help in understanding how to balance them effectively in practical genomic applications.

**3.2.1 Accuracy vs. Interpretability**

One of the primary trade-offs observed involves balancing model accuracy with interpretability. Naive Bayes and Logistic Regression both demonstrated high levels of accuracy, with Naive Bayes being particularly outstanding. While Logistic Regression provides relatively clear interpretations of its outputs through the coefficients associated with each feature, Naive Bayes, despite its probabilistic transparency, can sometimes mask how interdependent features influence predictions due to its assumption of feature independence. In contrast, Decision Trees offer excellent interpretability by allowing researchers to trace decision-making paths clearly. However, this model exhibited lower accuracy and is susceptible to overfitting, potentially compromising the reliability of its interpretations unless properly constrained and tuned. Models like Random Forest and XGBoost, which are ensemble methods, typically provide higher accuracy but at the expense of reduced interpretability. The aggregation of numerous decision trees in these models obscures the specific contributions of individual features to the overall predictions.

**3.2.2 Computational Efficiency vs. Model Complexity**

Another critical trade-off involves computational efficiency against the complexity of the model, which affects training and prediction speeds as well as practical deployment capabilities. Naive Bayes is highly efficient, capable of processing large datasets quickly due to its straightforward computational approach that avoids iterative parameter adjustments. This efficiency makes it particularly suitable for scenarios where speed is crucial. On the other hand, models like Gradient Boosting and XGBoost require more complex calculations and iterative adjustments, significantly extending model training and prediction times. Nonetheless, their ability to minimize errors and optimize performance might justify the increased computational costs in contexts where superior predictive accuracy is vital. Similarly, SVM can be computationally demanding, especially with large

and complex multi-class datasets. Its performance heavily relies on the choice of kernel and parameter tuning, which can further increase computational demands.

### 3.2.3 Precision-Recall Balance

The balance between precision and recall is especially critical in medical and biological applications, where the consequences of false positives and false negatives are significant. Naive Bayes and Logistic Regression exhibited strong performances in both precision and recall, making them suitable for applications where errors in either direction are costly. Models like Random Forest and XGBoost, while robust overall, may need tuning to enhance either precision or recall, depending on the application's specific needs. Their ensemble nature generally reduces variance, which can improve recall but sometimes at the cost of precision. Models such as SVM and Decision Trees struggled with recall, suggesting potential issues in detecting all positive cases within a dataset. This is a crucial consideration in areas like disease screening, where failing to identify a positive case can have severe implications. This trade-off analysis underlines the importance of the study's context in selecting machine learning models for genomic research. While Naive Bayes and Logistic Regression provide a good mix of efficiency and accuracy, their effectiveness may vary depending on the dataset characteristics and the study's specific demands. Conversely, models like Random Forest and XGBoost may be preferred where accuracy is paramount, though they may require compromises in terms of interpretability and computational efficiency.

### 3.3 Challenges and Limitations

While the study highlighted the efficiency and robustness of several machine learning models in classifying genomic sequences, it also revealed significant challenges and limitations inherent to some models, particularly Decision Trees and Support Vector Machines (SVMs). These challenges underscore the complexities involved in adapting traditional machine learning models to the nuanced demands of genomic data. Decision Trees, known for their intuitive decision-making paths, encountered issues of overfitting when applied to the genomic datasets. Overfitting occurs when a model learns the details and noise in the training data to an extent that it negatively impacts the performance of the model on new data. This is particularly problematic in genomic applications where the objective is to generalize findings across different, yet related, genomic sequences.

The model's tendency to overfit might be due to its capacity to create highly complex trees that perfectly describe the training data but fail to predict future observations accurately. While this can often be mitigated by setting constraints on the model parameters, such as limiting the depth of the tree or requiring a minimum number of samples at a leaf node, these adjustments may not always be sufficient to prevent overfitting entirely. Ensemble methods like Random Forest, which integrate multiple Decision Trees to make a final prediction, help in addressing this challenge by averaging the results of individual trees, thus reducing the risk of overfitting. This technique not only stabilizes the predictions but also improves the robustness of the model by aggregating a diverse set of decision paths and outcomes.

The SVM model exhibited lower performance, which can be attributed to several factors. One major challenge is selecting the appropriate kernel type and tuning its parameters, which are critical in defining the decision boundaries between different classes. Genomic data, characterized by its high dimensionality and complex pattern structures, requires a carefully chosen kernel to capture the essential relationships in the data effectively. The linear kernel may be too simple to model the complex relationships between features in genomic sequences, while more sophisticated kernels, like the radial basis function (RBF), require careful tuning of their parameters to avoid overfitting. The choice of kernel and its settings significantly influence the SVM's ability to generalize from training data to unseen data, a critical aspect in genomic studies where the ultimate goal is to apply findings across different biological contexts. Moreover, SVMs are also known for their computational intensity, especially when dealing with large datasets and multi-class classification scenarios. This can lead to longer training times and increased computational costs, which might not be feasible in all research settings.

### 3.4 Strategic Recommendations and Future Directions

The exploration and comparison of various machine learning models for genomic sequence classification have illuminated their potential and exposed their limitations. In order to enhance the effectiveness and reliability of these models in genomic research, a series of strategic recommendations are proposed. These recommendations focus on addressing the current challenges faced by the models, capitalizing on their strengths, and paving the way for innovative approaches in genomic analysis. Improving the robustness of machine learning models in genomic applications is crucial for their successful implementation. One approach to achieve this is through the adoption of advanced regularization strategies. These strategies are particularly vital for models like Decision Trees and SVM, which are prone to overfitting. By implementing regularization, the complexity of these models can be constrained, enhancing their generalizability and enabling them to perform more reliably on new, unseen genomic data.

Another approach to enhancing model robustness is through hybrid modeling techniques. These techniques involve combining the predictive power of multiple models to achieve better overall performance. For instance, integrating ensemble methods with high-precision models such as Naive Bayes could balance the

accuracy and robustness of the predictions. This could be executed by stacking various models where the outputs of some serve as inputs to others, creating a potent and comprehensive predictor. The computational demands of models like SVM and ensemble methods are significant, especially when handling large genomic datasets. To address this, there is a pressing need to enhance computational efficiency. One solution is the development of more efficient algorithms. This could involve optimizing the implementation of kernels in SVMs or improving the tree construction process in Random Forests, thereby reducing both computational time and resource consumption.

Furthermore, leveraging modern parallel computing technologies can significantly expedite data processing and model training. By implementing machine learning algorithms capable of parallelization, the runtime can be significantly reduced, making the application of these complex models more viable in large-scale genomic studies. These strategic recommendations, if implemented, could profoundly impact the field of genomic research by enhancing the performance, reliability, and applicability of machine learning models. By addressing these key areas, researchers can harness the full potential of machine learning to uncover new insights in genomics and push the boundaries of what is currently possible in genetic analysis and personalized medicine. These efforts will not only refine the predictive capabilities of these models but also broaden their applicability, leading to advancements that could transform our understanding and treatment of genetic disorders.

# 4. CONCLUSION

This study comprehensively evaluated several machine learning models for classifying human DNA sequences, highlighting each model's strengths and limitations. Through rigorous testing, models like Naive Bayes, Logistic Regression, Random Forest, and XGBoost demonstrated high accuracy and robustness, while Decision Trees and SVM faced challenges with overfitting and computational efficiency. Addressing these challenges with advanced regularization, optimized algorithms, and parallel computing is crucial for enhancing model performance in genomic research. The integration of machine learning in genomics promises improvements in genomic data classification and our understanding of genetic functions. This advancement is vital for personalized medicine, leading to better diagnostics, targeted treatments, and improved patient outcomes. Future progress in genomic research will depend on collaboration between bioinformatics, machine learning, and genetics to harness these technologies' full potential, enhancing genomic classification accuracy and efficiency and contributing to our understanding of complex biological systems and human health.

# REFERENCES

[1] P. Tolani, S. Gupta, K. Yadav, S. Aggarwal, and A. K. Yadav, "Big data, integrative omics and network biology," *Adv. Protein Chem. Struct. Biol.*, vol. 127, pp. 127–160, 2021.
[2] D. S. Bailey and G. I. Johnston, "Impact of genomics on the discovery and development of modern medicines," in *Genetics of Common Diseases*, Garland Science, 2020, pp. 241–261.
[3] A.-F. A. Mentis and L. Liu, "Global impact and application of Precision Healthcare," in *The New Era of Precision Medicine*, Elsevier, 2024, pp. 209–228.
[4] U. Radzikowska *et al.*, "Omics technologies in allergy and asthma research: An EAACI position paper," *Allergy*, vol. 77, no. 10, pp. 2888–2908, 2022.
[5] H. Satam *et al.*, "Next-generation sequencing technology: current trends and advancements," *Biology (Basel).*, vol. 12, no. 7, p. 997, 2023.
[6] L. Bai, Y. Wu, G. Li, W. Zhang, H. Zhang, and J. Su, "AI-enabled organoids: Construction, analysis, and application," *Bioact. Mater.*, vol. 31, pp. 525–548, 2024.
[7] P. Crovari *et al.*, "GeCoAgent: a conversational agent for empowering genomic data extraction and analysis," *ACM Trans. Comput. Healthc.*, vol. 3, no. 1, pp. 1–29, 2021.
[8] A. Sharma and R. Kumar, "Recent Advancement and Challenges in Deep Learning, Big Data in Bioinformatics," in *Blockchain and Deep Learning: Future Trends and Enabling Technologies*, Springer, 2022, pp. 251–284.
[9] M. K. Gupta *et al.*, "Sequence Alignment," *Bioinforma. Rice Res. Theor. Tech.*, pp. 129–162, 2021.
[10] J. K. Chaudhari, S. Pant, R. Jha, R. K. Pathak, and D. B. Singh, "Biological big-data sources, problems of storage, computational issues, and applications: a comprehensive review," *Knowl. Inf. Syst.*, pp. 1–51, 2024.
[11] U. Ullah and B. Garcia-Zapirain, "Quantum Machine Learning Revolution in Healthcare: A Systematic Review of Emerging Perspectives and Applications," *IEEE Access*, 2024.
[12] Y. Cao, T. A. Geddes, J. Y. H. Yang, and P. Yang, "Ensemble deep learning in bioinformatics," *Nat. Mach. Intell.*, vol. 2, no. 9, pp. 500–508, 2020.
[13] N. S. Kiran, C. Yashaswini, R. Maheshwari, S. Bhattacharya, and B. G. Prajapati, "Advances in Precision Medicine Approaches for Colorectal Cancer: From Molecular Profiling to Targeted Therapies," *ACS Pharmacol. \& Transl. Sci.*, vol. 7, no. 4, pp. 967–990, 2024.
[14] S. Maleki Varnosfaderani and M. Forouzanfar, "The Role of AI in Hospitals and Clinics: Transforming Healthcare in the 21st Century," *Bioengineering*, vol. 11, no. 4, p. 337, 2024.
[15] J. F. Uleman, R. Quax, R. J. F. Melis, A. G. Hoekstra, and M. G. M. O. Rikkert, "The need for systems thinking to advance Alzheimer's disease research," *Psychiatry Res.*, vol. 333, p. 115741, 2024.
[16] V. Gambardella *et al.*, "Personalized medicine: recent progress in cancer therapy," *Cancers (Basel).*, vol. 12, no. 4, p. 1009, 2020.

[17] J. Peng, E. C. Jury, P. Dönnes, and C. Ciurtin, "Machine learning techniques for personalised medicine approaches in immune-mediated chronic inflammatory diseases: applications and challenges," *Front. Pharmacol.*, vol. 12, p. 720694, 2021.

[18] K. Huang, C. Xiao, L. M. Glass, C. W. Critchlow, G. Gibson, and J. Sun, "Machine learning applications for therapeutic tasks with genomics data," *Patterns*, vol. 2, no. 10, 2021.

[19] M. Barshai, E. Tripto, and Y. Orenstein, "Identifying regulatory elements via deep learning," *Annu. Rev. Biomed. Data Sci.*, vol. 3, pp. 315–338, 2020.

[20] A. A. Joshi and R. M. Aziz, "Deep learning approach for brain tumor classification using metaheuristic optimization with gene expression data," *Int. J. Imaging Syst. Technol.*, vol. 34, no. 2, p. e23007, 2024.

[21] N. Vasani, "Human DNA Data." 2022.