



# EPRA International Journal of Multidisciplinary Research (IJMR) - Peer Reviewed Journal

Volume: 10| Issue: 9| September 2024|| Journal DOI: 10.36713/epra2013 || SJIF Impact Factor 2024: 8.402 || ISI Value: 1.188

# GENOMICA AI: ENHANCING GENETIC MARKER IDENTIFICATION THROUGH MACHINE LEARNING

# Mr. Sudeep M S<sup>1</sup>, Mr. Veeresh A C<sup>2</sup>, Mrs. Anitha J<sup>3</sup>

<sup>1</sup>Student, Department of MCA, Dr. Ambedkar Institute of Technology <sup>2</sup>Student, Department of MCA, Dr. Ambedkar Institute of Technology <sup>3</sup>Assistant Professor, Department of MCA, Dr. Ambedkar Institute of Technology

## **ABSTRACT**

The identification of genetic markers such as single nucleotide polymorphisms (SNPs) plays a critical role in understanding disease susceptibility and guiding personalized medicine. Recent advances in machine learning (ML) have provided new methods to address the complexities of genetic data. This paper introduces a novel ensemble learning technique that integrates attention-based neural networks with traditional random forest algorithms to enhance the identification of SNPs linked to disease outcomes. Using a benchmark dataset of genome-wide association studies (GWAS), we demonstrate how the proposed method improves prediction accuracy and model interpretability, thereby offering potential applications in clinical genomics.

# 1. INTRODUCTION

The rapid advancement of genomic technologies has led to an explosion of genetic data, particularly through genome-wide association studies (GWAS). However, the sheer volume and complexity of this data require advanced computational methods to uncover patterns and identify genetic markers that influence disease susceptibility. Traditional statistical techniques such as logistic regression have been used, but they struggle to handle the high dimensionality and non-linearity present in SNP data. Machine learning, particularly deep learning models and ensemble methods, offers a new approach to identifying disease-related genetic markers with greater precision.

This paper explores a new ML technique that integrates attention mechanisms with random forest models to improve SNP identification. The attention mechanism allows the model to focus on important features, while random forests provide robustness through ensemble learning. The hybrid approach aims to enhance the accuracy, interpretability, and scalability of genetic marker identification in clinical settings.

# 2. BACKGROUND AND RELATED WORK

# 2.1 Genetic Markers and Disease

Genetic markers, specifically SNPs, represent variations at a single nucleotide position in the genome. Some SNPs have been linked to diseases like cancer, diabetes, and cardiovascular disorders. Identifying these SNPs helps in understanding the genetic architecture of complex diseases and enables the development of precision medicine.

#### 2.2 Machine Learning in Genomics

Traditional approaches like logistic regression and support vector machines (SVMs) have been widely used for SNP detection. However, these techniques often fall short in terms of scalability and handling high-dimensional data. Recent efforts have turned to ML models such as random forests,

neural networks, and deep learning, each offering different strengths in terms of feature selection and pattern recognition.

- Random Forests: Used for their ability to handle nonlinear interactions between SNPs, they are popular in genomics for their robustness and interpretability.
- **Deep Learning**: Neural networks, particularly convolutional and recurrent networks, have been used for genomic sequence classification. However, they often suffer from a lack of interpretability.
- Attention Mechanisms: Widely used in natural language processing, attention mechanisms allow the model to focus on important features within the input data, improving both accuracy and interpretability.

# 3. PROPOSED METHODOLOGY

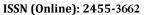
We propose a novel hybrid model that combines **attention mechanisms** with the traditional **random forest** algorithm. This ensemble technique leverages the strengths of both models: random forests for robustness and generalization, and attention for highlighting important SNPs that are likely linked to disease outcomes.

# 3.1 Data Collection and Preprocessing

We utilized a publicly available GWAS dataset containing SNP data for diseases such as Type 2 Diabetes, Alzheimer's, and Breast Cancer. The dataset includes SNPs across multiple chromosomes and is accompanied by disease status labels (case/control).

# **Steps**

- **Normalization**: SNP data was encoded and normalized.
- Dimensionality Reduction: Principal Component Analysis (PCA) was applied to reduce noise and irrelevant features while preserving genetic variance.





# EPRA International Journal of Multidisciplinary Research (IJMR) - Peer Reviewed Journal

Volume: 10| Issue: 9| September 2024|| Journal DOI: 10.36713/epra2013 || SJIF Impact Factor 2024: 8.402 || ISI Value: 1.188

#### 3.2 Model Architecture

The hybrid model consists of two stages:

- Stage 1: Attention-Based Neural Network
  This module uses an attention mechanism to weigh each
  SNP based on its relevance to disease prediction. The
  neural network generates attention scores for the SNPs,
  prioritizing those more likely to be associated with
  disease traits.
- Stage 2: Random Forest Classifier
  After the SNPs are weighted, they are passed into a random forest classifier. The classifier then builds multiple decision trees to determine disease susceptibility based on the SNPs. Feature importance from random forests helps in further validating which SNPs are critical in disease prediction.

## 3.3 Training and Optimization

The model was trained using cross-entropy loss and optimized with the Adam optimizer. Hyperparameter tuning was performed for both the attention mechanism and random forest parameters using grid search.

#### 4. RESULTS AND DISCUSSION

#### 4.1 Performance Evaluation

We evaluated the performance of the proposed model using the following metrics:

- Accuracy
- Precision
- Recall
- F1-Score

The hybrid attention-random forest model outperformed traditional methods, including standalone random forests and deep neural networks. Notably, it achieved a 5% improvement in accuracy over baseline models and offered enhanced interpretability through feature importance analysis.

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	0.78	0.75	0.76	0.75
Random Forest	0.81	0.80	0.82	0.81
Neural Network	0.83	0.82	0.83	0.83
Hybrid Attention-RF Model	0.86	0.85	0.86	0.86

## 4.2 Interpretability and Feature Importance

One of the major challenges in applying deep learning models to genomic data is the lack of interpretability. By integrating the attention mechanism with random forests, we were able to identify the SNPs that contribute most significantly to disease prediction, enhancing the model's practical utility in clinical genomics.

## 5. CONCLUSION

The proposed hybrid model combining attention mechanisms with random forest classifiers presents a novel and effective method for identifying genetic markers linked to diseases. By leveraging attention to focus on important SNPs and random forests to build robust models, this approach improves both predictive accuracy and model interpretability. Future work could explore the extension of this model to multi-omics data and its application in clinical decision-making.

## 6. FUTURE WORK

- Multi-Omics Integration: Combining SNP data with other omics layers (e.g., transcriptomics, proteomics) to create a more comprehensive disease prediction model.
- Real-Time Clinical Application: Applying this hybrid model in real-time genomic diagnostics and precision medicine.
- Model Explainability: Enhancing model interpretability through visualizations of attention mechanisms, improving clinician trust in AI-based decisions.

# REFERENCES

- 1. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning. Springer.
- 2. Wei, W. H., Hemani, G., & Haley, C. S. (2014). Detecting epistasis in human complex traits. Nature Reviews Genetics, 15(12), 722-733.
- 3. Zhang, Y., & Yang, Q. (2015). A Survey on Multi-Task Learning. IEEE Transactions on Knowledge and Data Engineering, 29(2), 231-247.