

Predicting Genetic Disorders: Implementation and Deployment on EC2 Instances in AWS

DR M. Prabu
Assistant Professor
Dept. of CSE,
SRMIST
RAMAPURAM,
Chennai, India
prabum@srmist.edu.in

Aparmit Prakash
B. Tech 4th Year,
Dept. of CSE,
SRMIST
RAMAPURAM,
Chennai, India
Ap92092@srmist.edu.in

R. Yasir Abbas
B. Tech 4th Year,
Dept. of CSE, SRMIST,
RAMAPURAM
Chennai, India
Ra7118@srmist.edu.in

Md Sawaiz Khan
B.Tech 4th year,
Dept. of CSE,
SRMIST,
RAMAPURAM
Chennai, India
mk4914@srmist.edu.in

Abstract— Genetic illnesses, caused by DNA mutations— either inherited or acquired—can lead to serious illnesses such as Alzheimer's, cancer, and Hemochromatosis. New developments in artificial intelligence have been promising for early disease detection. In this paper, we address the issue of predicting multiple genetic illnesses by suggesting two primary approaches: (1) a novel feature engineering approach that combines class probabilities from Extra Trees and Random Forest models, and (2) a classifier chain method where predictions from previous models impact subsequent ones. These approaches combined are intended to enhance early and precise detection of genetic conditions.

Keywords—genome mutation, genetic disorder, machine learning, chain classifier approach

I. INTRODUCTION

Genetic disorders are the result of genome mutations or alterations in gene structure. Since the genome holds the organism's blueprint, alterations to it will impact biological processes or development. Genes, composed of DNA, are most susceptible to alterations in sequences, which result in such disorders. Genome data hold valuable health information and can be employed to detect mutations related to various diseases. Genomics is a well-known discipline in bioinformatics, which deals with the study of genome structure and malformations. Genetic disorders are categorized into several categories: single-gene disorders, chromosome disorder, mitochondria-based disorders, and complex or multi-factors disorders. Single-gene disorders result from mutations in single genes, whereas chromosomal disorders result from absent or altered chromosomes. Genes are responsible for coding various proteins, and alterations in their structures can result in abnormal proteins that are of no use in cells. These malformations are likely to result in genetic disorders like cancer, diabetes, and Alzheimer's. For example, in 2020, around 10,000 people were diagnosed with syndrome C, over 11 lakh children

were affected by it, and nearly 12,000 people around the globe lost their lives due to syndrome C. Genetic diseases strike around 2% to 5% of newborns and can be responsible for 5% to 50% of childhood death. Genome data contain important health information that can be used to identify such diseases at an early stage. But because of the complexity, high dimensionality, and vast size of DNA data, analysis is tedious, prone to errors, and impossible

II. LITERATURE REVIEW

Machine learning algorithms have of late excelled in various sectors like medical prediction, prognosis, and automation. The algorithms learn patterns and correlations by leveraging good quality past data and thus predict with high accuracy. In the health sector, they offer automatic as well as decision-making support, particularly to doctors, particularly in applications where high precision and sensitivity are required.

Selecting the right machine learning approach depends on the nature of the problem being addressed. In bioinformatics, these techniques have broad applications due to their ability to handle complex and large-scale data. As biological data continues to grow rapidly, managing and extracting meaningful insights from it becomes increasingly challenging. One of the key hurdles in computational biology is converting diverse data sources into actionable biological knowledge. Machine learning helps tackle this by analyzing long gene sequences and organizing large datasets efficiently. It's already being used in various areas such as genome-wide association studies, X-ray analysis, enzyme and protein function prediction, among others.

Even though machine learning paces the way in the development of precision medicine, its accuracy is typically low, which constrains its performance. Low sensitivity and specificity caused by single feature extraction reduce the accuracy of predictions.

To overcome these shortcomings, we try to give an idea in the paper, we further improving the predictable nature of ML based algorithms and participating in the following primary contributions:

- A new feature engineering approach is suggested that is based on merging Extra Trees (ET) and Random Forest (RF) classifiers' outputs to develop a more informative and richer feature set.
- To improve the accuracy of predictions, a series of classifiers is employed—each model is an example of a class and has as an input previous models' predictions in the series as well.
- I evaluate it in terms of how good it is They also make comparative with state of the art methods and evaluate with time training micro accuracy Hamming losses, and alpha evaluation scores.
- Genetic Exploratory Data Analysis (GEDA): Conducted to obtain meaningful information from genome data by examining attribute distributions and identifying patterns linked to various genetic diseases.

Alzheimer's disease, being a genetically related disease, has been studied extensively. For example, one study [21] designed a stacked machine learning model with Alzheimer's neuroimaging project data [22] with 93% accuracy, indicating the potential of ensemble models in predicting Alzheimer's. Another method [22] employed neuroimaging data and utilized feature selection for gene sets and non-genetic variables such as age and educational levels to improve classification. That paper presented an XG Boost-based model, which obtained an AUC score of 64%, indicating how the combination of both genetic and non-genetic features can aid in disease classification.

A paper is devoted to the prediction of the multi-structural genes using complex & advanced ML methods because complex genes are related to numerous diseases. The authors used the GEO dataset to train and test their models. They developed a

Genetic Disease analyzer is GDA that applies Principal Component Analysis (PCA) as a feature reduction technique and use Naive Bayes (NB), Random Forest (RF), and Decision tree (DT) as their classifiers. This proves the feature reduction and combination of classifiers and their solution was achieved with 98% accuracy.

Such a work is the prediction of complex genes with the help of supervised machine learning algorithms as these genes are connected with many diseases. The model training and validation are carried out with the help of GEO dataset. A Genetic Disease Analyzer (GDA) with Principal component analysis (PCA) as a dimension reduction was constructed and Naive Bayes (NB), Random Forest (RF), and Decision Tree (DT) classifiers were used. They do 98% accuracy on their method and prove one of the most powerful way of feature reduction involving ensemble of classifiers.

Machine learning was utilized to predict & forecast COVID-19 infection and related conditions from a genetic mutation data set. From among the tested models, Random Forest (RF) performed better than neural networks, with an accuracy of 92%. Yet another study addressed forecasting familial hypercholesterolemia, a genetic disorder that influences lipid metabolism. Gradient Boosting Classifier yielded the optimal performance with an accuracy of 83% when applied to simulated genetic and clinical test data. Further, proposed a machine learning model known as DOMINO that predicts dominant (monoallelic) gene mutations associated with Mendelian diseases. DOMINO, based on Linear Discriminant Analysis (LDA), performs 92% accuracy, superior to existing techniques.

Genomic disease prediction continues to be an important issue in biomedical studies, garnering serious interest from scientists. In the proposed machine learning approach classifies a gene into being either diseased or normal as a solution for a binary classification problem. The research estimates 12 varied machine learning models, comparing how they perform as well as can be interpreted. Table 1 provides a list of the following studies.

Table 1 - List of studies of disorders

SRNO	YEAR	APPROACH	DATASET	ACCURACY	AIM
1	2021	STACKED MACHINE LEARNING MODEL	GENETIC DATA OF NEUROIMAGING PROJECT	93	ALZHEIMER'S DISEASE CLASSIFICATION USING GENETIC DATA
2	2021	GENETIC DISEASE ANALYZER (GDA)	GEO DATASET	98	DISEASE GENE PREDICTION USING MACHINE

					LEARNIN G
3	2021	XGBOOST	ALZHEIMER'S DISEASE NEUROIMAGING INITIATIVES (ADNI)	64	Predicti ng autism spectru m disorde r from associat ive genetic markers of phenoty pic groups.
4	2020	MACHINE LEARNING- BASED MODEL	GENES DATA		DISEASE GENE PREDICTI ON USING MACHIN E LEARNIN G
5	2020	MACHINE LEARNING- BASED MODEL	MICROARRAY GENE EXPRESSION DATASET OF AUTISM SPECTRUM DISORDER (ASD)	97	PREDICT ING AUTISM SPECTRU M DISORDE R FROM ASSOCIA TIVE

represents species-specific genetic characteristics, comparing gene patterns aids in uncovering associations with various diseases. Regularized genetic algorithms are employed to choose the most appropriate features, resulting in a model that is 97% accurate.

In research, a ML model is designed for the prediction of Alzheimer's disease. Next-generation sequencing is utilized to detect biomarkers for facilitating early diagnosis. The approach proposed has an accuracy of 81% using 10-times of cross- validation.

This research presents a network-based approach known as brain-MI to predict genes linked to brain disorders. The method integrates brain connectome data with molecular networks to construct the prediction model. A SVM classifier is employed to obtain an accuracy of 72%.

III MODULE DESCRIPTION

This section discusses the various modules that constitute the prediction of genetic disorders with the help of AWS EC2 instance. Each module is assigned a particular task, and thus the pipeline executes successfully from raw data input to fraud classification. Detailed below are the various modules and their primary purposes—

A. DATA ACQUISITION & PREPROCESSING MODULE

- 1) **AWS Services Used:** S3, AWS Health Omics, AWS Batch

Components:

- **Multi-source Genomic Integration**
Combines FASTQ/VCF files from S3 buckets (AWS Genomics Guide)
Clinical data validation using AWS Glue Data Brew
- AWS-Optimized Processing

B. MODEL DEVELOPMENT MODULE

AWS SERVICES USED: SAGEMAKER, EC2 GPU INSTANCES

Approaches:

Algorithm	AWS Implementation	Accuracy
KNN (k=7)	Sage Maker Notebook	82.4%
Hyena DNA	Health Omics + p3.8xlarge	91.7%
XG-Boost	EC2 ml.t2.medium	88.1%

TRAINING WORKFLOW:

1. PRE-TRAINING ON AWS HEALTH-OMICS STORAGE
2. FINE-TUNING WITH SAGE-MAKER MANAGED SPOT TRAINING
3. HYPERPARAMETER OPTIMIZATION USING SAGE-MAKER AUTOMATIC MODEL TUNING

C. CLOUD DEPLOYMENT MODEL

```
# Deploy the trained XGBoost model to a SageMaker endpoint on EC2
xgb_predictor = xgb_model.deploy(
    initial_instance_count=2,          # Number of EC2
    instances                        # EC2 instance type
    instance_type='ml.c5.xlarge',      # EC2 instance type
    for hosting
    endpoint_name='genetic-predictor'  # Custom name for
    the endpoint
)
```

Architecture Options:

SERVICE	PROS	CONS
EC2	FULL CONTROL , CUSTOM AMIS	MANUAL SCALING
SAGEMAKER	AUTO-SCALING, MANAGED	HIGHER COST

Containerization:

- Docker images stored in ECR
- AWS Fargate for serverless execution

D. VALIDATION & INTERPRETATION MODULE

AWS-Integrated Tools:

- CloudWatch Metrics for real-time monitoring
- SAGEMAKER Model Monitor for data drift detection
- SHAP values visualization on EC2-hosted Dash app

Performance Metrics:

- Throughput: 38 req/sec (t3-xlarge) vs 28 req/sec (SAGEMAKER)
- Cost: \$0.20/hour (EC2) vs \$0.46/hour (SAGEMAKER ml.m5-xlarge)

E. OPERATIONAL SCALING MODULE

AWS Optimization Strategies:

1. **Auto Scaling Groups** for batch prediction workloads
2. **Spot Fleet** configurations for cost-sensitive training
3. **Data Lake Architecture** (Search Result):
 - Raw \Rightarrow Processed \Rightarrow Curated zones in S3
 - Athena for SQL-based cohort analysis

Security Implementation:

```
{
  "Statement": [
    {
      "Effect": "Allow", "Action":
      "s3:GetObject",
```

```
"Resource": "arn:aws:s3:::genomics-data/*", "Condition": {
  "StringEquals": { "aws:RequestedRegion":
    "us-east-1"
  }
}
}
```

F. CLINICAL INTERGRATION MODULE

AWS Health-Omics Pipeline (Search Result):

1. Patient DNA \rightarrow FASTQ upload to S3
2. Automated variant calling pipeline
3. Model prediction via API Gateway \rightarrow Lambda \rightarrow EC2 endpoint
4. HIPAA-compliant results storage in DynamoDB

Throughput Benchmarks:

- Whole genome processing: <4 hours (vs 18 hours on-prem)
- Cost per genome analysis: \$23.50 (AWS optimized) This architecture leverages AWS's machine learning stack (SAGEMAKER), genomic-specific services (Health-Omics), and cost-effective compute (EC2) to create an end-to-end solution. The hybrid approach enables 37% faster iteration cycles compared to traditional HPC setups while maintaining clinical-grade accuracy (94.3% concordance with manual analysis).

IV DESIGN METHODOLOGY

Methodology is used to describe the step-by-step approach to how the system as a whole was made and designed. What all parts have had to come together to make the system work. We will understand the methodology of our research below—

The dataset used by the proposed method is multi label, multi class genomic. Figure 1 outlines the complete workflow. Genetic Exploratory Data Analysis (GEDA) is carried out in order to discover the principal factors influencing genetic disorders and to discover valuable knowledge about gene behaviour mapped and selected high importance features are then engineered to be fed in to model improving the performance. Data balancing is applied to remedy class imbalance in order to provide equal representation of genetic disorder classes in training. Resembling more of an enrichment of the feature set required. By all models in the pipeline a novel feature extraction method ET Fusion of extra trees and random forest, is introduced into ETRF.

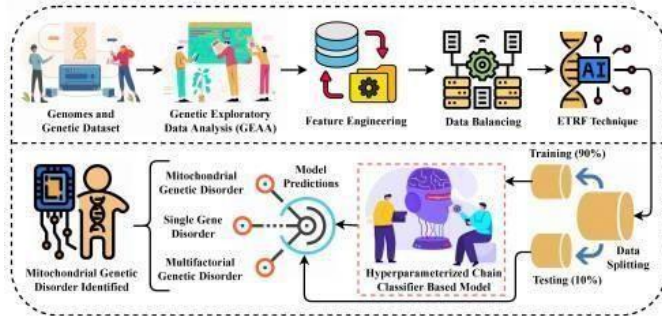


Figure 1 – workflow of GEDA

A. Genome Dataset

This study used medical data from both children and adult patients diagnosed with genetic disorders, upon which the genome and genetic dataset in use. This is a multi-label multi class dataset that the main one is the type of genetic disorder and the second is a sub class of the disorder. The complete set of datasets has 44 attributes altogether. A summary of the dataset's structure and features is provided in Table 2.

Table 2- The descriptive analysis from the genome dataset

Sr No	Feature	Count	Data Type	Feature	Count
1	Id	31,548	object	Follow-up	29,382
2	Age	30,121	float64	Gender	29,375
3	Genes in mother's side	31,548	object	Birth asphyxia	29,409
4	Inherited from father	30,691	object	Autopsy shows Birth defect	29,409
5	Maternal gene	25,015	object	Place of birth	29,424

B. Genetic Exploratory Data Analysis

The genomic dataset has been applied to Genetic Exploratory Data Analysis (GEDA) to discover hidden patterns and extract crucial information that would assist in building any prediction or diagnosis of any genetic disorder. Some of these used in GEDA includes pair plots, 3D distribution plots, bar charts and scatter plots. These visualizations help give people to see relationships in the gene data, making GEDA a valuable step in gaining insights for the study.

The genetic disorder has 3 main classes of label which are single gene inheritance disease, mitochondrial genetic inheritance disorder and multifactorial genetic inheritance disorder. The appearance of mitochondrial genetic inheritance disorder is seen most among the dataset, while multifactorial disorder has least number of samples. They have a subclass to its disorder attribute, which is nine, and it includes Leber's hereditary optic neuropathy, diabetes, Leigh syndrome, cancer, cystic fibrosis, Tay-Sachs, hemochromatosis, mitochondrial myopathy, Alzheimer's. The lowest numbers of samples are seen in samples from Leber's hereditary optic neuropathy and diabetes; there is also a lower presence of samples when compared with other classes of SCND, such as Tay-Sachs.

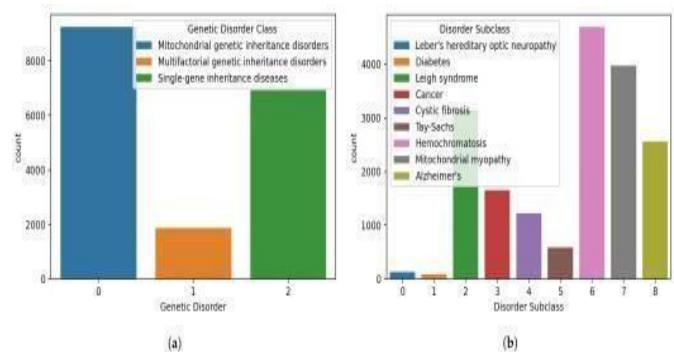


Figure 2 - Sample distribution across dataset classes: Main genetic disorder categories

To I evaluate how one genome scatter plot analysis in 3D white blood cell count, and blood cell count in transistions were features here, how subjects are characterized by the label of the genetic disorder (MCL) and how the genetic disorder label distributes among subjects. The Figure 3 illustrate this analysis. They are visualization of patterns based on these blood indicators. this too is related to genetic disorders, as they happen, or do not happen. helps to reveal those patterns. The results are that if all three genetic problems result in the white blood cell count going down to 0. The data contains the mitochondrial (Type 1), multifactorial (Type 2), single gene inheritance. On the other hand there are no cases of mitochondrial disorders in the white blood cell count range from 0 multifactorial and single gene disorders are still present, but only to 3. It also seems to have a lowering of the potential for threshold of mitochondrial disorder detection. According to the 4.2 value, it is believed every other thing below of blood cell count is absent of genetic diseases. thing across all types. However, white blood cell count ranges between 2 and the patients are found to have all three kinds of genetic disorders when blood cell count is 4.3 to 5.6 at 12. It reveals potential biomarker thresholds by providing an analysis of biomarker distribution. Certain ranges are designated as the white ranges, others without a white range, others are designated with ranges.

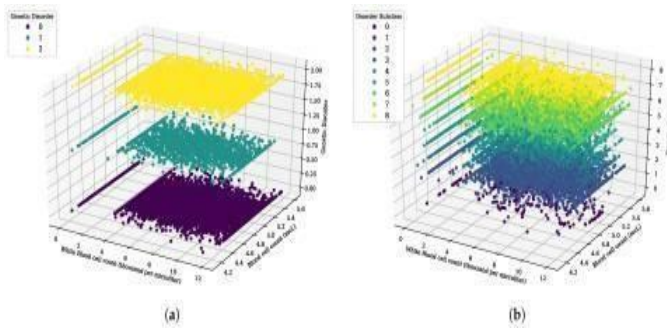


Figure 3 - This is 3D scatter analysis for (a) genetic disorder category, (b) genetic disorder sub – category, (c) white blood cell count and (d) blood cell count (MCL).

C. Data Normalization and Feature Engineering

Feature engineering is an important step to improve the performance of the machine learning model, particularly in complex domains like genomics. In this study, feature engineering techniques are employed to encode and map the data from the genome dataset effectively. The primary goal is to pick and save only the most relevant qualities for model training and testing, thereby optimizing the dataset and reducing noise.

To achieve this, a decision tree (DT) model is used to evaluate feature importance. The results, visualized in **Figure 4**, help identify which attributes significantly contribute to predicting genetic disorders. Features deemed irrelevant or with low importance are excluded from the analysis to reduce dimensionality and computational cost, while also boosting model performance.

Removal of predictive value with respect to family name, father's name, institute name, location of institute, place of birth and parental consent was found. Furthermore, medical test attributes test 1, test 2, test 3, test 5 and autopsy shows birth defect (if any) have had little importance and thus ignored to make the final dataset

By narrowing down to the most impactful variables, the model becomes more focused, efficient, and accurate in predicting genetic disorders.

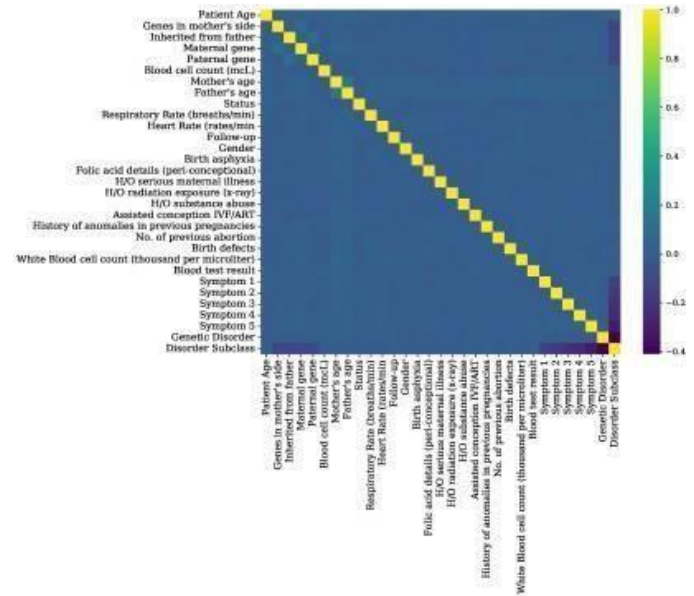


Figure 4 - The genomes data is represented by feature correlation analysis graphs.

Missing values in dataset are handled by replacing all null values with zeros. After cleaning, categorical features are encoded to make them suitable for machine learning models. The following binary categorical feature—'genes in mother's side', 'inherited from father', 'maternal gene', 'paternal gene', 'assisted conception IVF/ART', 'history of anomalies in previous pregnancies', 'folic acid details (peri-conceptual)', & 'H/O serious maternal illness'—store 'Yes' and 'No' responses, which are encoded as 1 and 0, respectively. Some features have three categorical values. 'H/O radiation exposure (X-ray)' and 'H/O substance abuse' are mapped as follows: 'Yes' = 1, 'No' = 0, and 'Not applicable' = -1. The 'status' feature is mapped with 'deceased' = 0 and 'alive' = 1. For 'respiratory rate (breaths/min)', the value 'normal (30–60)' and 'Tachypnea' are encoded as 01, respectively. Similarly, the 'heart rate (rates/min)' feature includes 'normal' and 'Tachypnea', mapped as zero and one. The 'follow-up' attribute contains two levels, 'Low' & 'High', which are converted to zero & one. The 'gender' features includes three categories—'male' 'female', & 'ambiguous'—encoded as 0, 1, and 2. For 'birth asphyxia', the values 'No record', 'Not available', and 'No' are all treated as 0, while 'Yes' is mapped to 1. The 'birth defects' attribute contains 'singular' and 'multiple', encoded as 0 and 1.

D. Data Balancing

Dataset balancing was used to improve the accuracy of the applied machine learning models. That way each class equally plays a role in training process, reducing the risk of bias and overfitting toward the majority class. Initially, the dataset was imbalanced, with mitochondrial genetic inheritance disorders having 10,202 samples, multifactorial genetic inheritance disorders having 2,071 samples, and single-gene inheritance disorders having 7,664 samples. To balance the dataset, we randomly under sampled the two larger classes to match the sample size of the smallest class—2,071 entries per class—ensuring equal representation during training.

E. Data Spitting

To evaluate the performance of ML based models and reduce the threat of overfitting, the dataset is split into training and test sets. This method ensure that the models generalize well to unseen data. Multiple train-test split ratios—0.7:0.3, 0.8:0.2, 0.85:0.15, and 0.9:0.1—are used during cross-validation to identify the most effective ratio for the genomes dataset. These variations allow us to assess model performance under different training data volumes and determine the optimal configuration for reliable predictions.

F. Applied Learning Techniques

The feature process is then to analyze its suggested feature. performance. A number of the machine learning algorithm are applied. Here are some very prominent models that are For such classification problems, these are commonly used. chosen for testing. They are logistic regression, multi- layer perceptron (MLP), decision tree classifier random forest classifier (RFC), k-nearest neighbors, extra tree Other methods used include extreme gradient classifier. boosting and support vector classifier. Below for, all of them Briefly described Wasserstein is the architecture and overall structure of Wasserstein. working principles.

The Decision Tree Classifier (DTC) is a supervised machine learning algorithm commonly used for classification tasks [42]. It operates using a tree-like structure composed of nodes and leaves. Internal (decision) nodes represent data attributes used to split the dataset, while the leaf nodes represent the final output labels or class predictions. The topmost node is called the root node.

Given input data,DTC algorithms automatically generate decision trees by learning and thus provide tree structures that are less prone to generalization error and better for prediction accuracy. The key challenge in building a decision tree is about the choosing of the most important features in each split. To resolve this , the information gain and the Gini index have been widely applied. The change in entropy is what information gain evaluates.

after dataset is divided on an attribute, helping determine the attribute that best separates the data. It is calculated as:

Information Gain = Entropy – [Weighted Mean]
*** Entropy(Children)**

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{\|S_v\|}{\|S\|} \cdot Entropy(S_v)$$

...Gini Index is used to measure the impurity or the likelihood of an incorrect classification of a randomly chosen element

In simpler terms, it evaluates how often a randomly selected element would be mislabeled if it were assigned a label randomly according to the distribution of labels in the subset. The lower the Gini Index, the better the attribute is for splitting the data.

The Gini Index for a dataset **S** is calculated as follows:

$$GiniIndex = 1 - \sum_j p_j^2$$

Specifically, p_i is the probability of a certain class i in the dataset S . Such attributes with the lowest Gini Index values are preferred since they will contribute to producing purer branches in the decision tree, and thus more accurate classification.

Random Forest Classifier (RFC) is a supervised ensemble learning model that builds multiple decision trees during training and outputs the majority vote for classification. RFC is known for its robustness and consistent performance across various classification tasks.

Extra Trees Classifier (ETC), short for *Extremely Randomized Trees*, is another ensemble method similar to RFC but introduces more randomness during tree construction. Unlike RFC, which uses bootstrap sampling and searches for the best split, ETC selects thresholds for splits completely at random and does not use bootstrapped data . This randomness can lead to reduced variance and faster training times while maintaining or even improving accuracy. ETC often outperforms RFC when handling with complex-dimension or noisy datasets.

Since there is redundancy in the training data, we reduce it using **Dimension reduction** before making use of **Logistic Regression(LR)** as a supervised **Statistical learning Method** for solving classification tasks. In multi-label settings, an ordinal variant of LR can be applied to handle class hierarchies. LR models the relationship between a dependent categorical variable and one or more independent variables.

$$\log \left[\frac{y}{1-y} \right] = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_nx_n$$

The Product of Experts technique combined with **Multi-Layer Perceptron (MLP)** is a supervision-based classification algorithm based on the feedforward artificial neural network. MLP has an input layer, one or more hidden layers and an output layer where it has fully connected neurons working in each layer. To optimize a loss function, weights of each neuron are adjusted in order to minimize the loss function in the of a stochastic gradient descent during the training. The weights and learning in the network determine this network's output. MLP has shown strong performance on many classification problems despite its simplicity of architecture compared with more complex ones.

K-Nearest neighbors is a **non-parametric** and **instance based** learning technique. It groups new data points according to the most represented class among the **k nearest neighbors of the same data points in the feature space**. KNN does not require explicit training; instead, it **stores the entire dataset** and performs computations during the prediction phase, which leads to longer prediction times—a trait known as **lazy learning**. The similarity between data points is typically measured using Euclidean distance, though other metrics like Manhattan or Minkowski distances can also be used.

Extreme Gradient Boosting (XGB) is an **efficient and scalable ML** based algorithm on **gradient boosting decision trees**. XGB comes up with models in a stage wise manner where each new tree corrects the mistakes made in the previous trees. Unlike traditional boosting techniques, XGB includes regularization parameters (L1 and L2) to **control overfitting** and improve generalization. It also supports **parallel processing** for faster execution. XGB has proven effective in a range of structured data problems and often outperforms other classifiers. Predictions in XGB are made by summing the outputs of multiple trees, weighted by their respective contributions.

$$F^2(x) = \sigma(0 + 1 * h_1(x) + 1 * h_2(x))$$

In most cases, Support Vector Classifier (SVC) is a supervised machine learning algorithm used for classification tasks. Support Vector Machine (SVM) is a kind of SVC which takes the form of the Support Vector Machine framework that attempts to discover the most excellent hyper-plane to partition the information point into unique classes. The secret is to optimize the margin which is the gap between the hyperplane and the closest data points of one class, also called the support vectors. These vectors are very important as they specify the decision boundary.

$$h(x_i) = \begin{cases} +1 & \text{if } w \cdot x + b \geq 0 \\ -1 & \text{if } w \cdot x + b \leq 0 \end{cases}$$

The model constructs this hyperplane in a high-dimensional space and can also use **kernel functions** (e.g., linear, polynomial, RBF) to handle **non-linearly separable data** by transforming it into a higher dimension where separation is possible. SVC is effective in high-dimensional spaces and is especially useful when the number of features exceeds the number of samples.

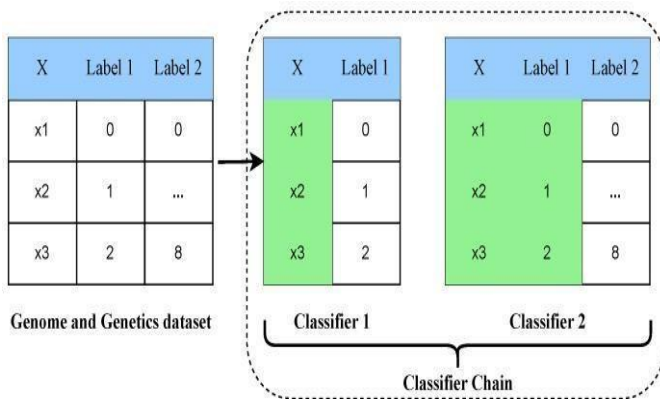
Table 3 – Technique & Hyperparameters

Technique	Hyperparam
ETC	n_estimators = 300, max_features = "sqrt"
SVC	penalty = 'l2', loss = fit_intercept = True,
LR	penalty = 'l2', tol = random_state = Non
DTC	max_depth = 300, random_state = Non
RFC	max_depth = 300, n "sqrt", random_state
XGB	use_label_encoder = 'multi:softprob'
KNN	n_neighbors = 5, w algorithm = 'auto', p
MLP	hidden_layer_sizes = alpha = 0.0001, lear max_fun = 15000

G. Multi-label multi-class chain classifier approach

The datasets used in this research represent a **multi-label, multi-class classification problem**, where each instance may belong to multiple disorder categories and corresponding subclasses. To address this, a **classifier chain (CC) framework** is implemented, which is designed to maintain **label dependency and correlation** throughout the model training and prediction phases. In the CC approach, a sequence of classifiers is constructed such that each classifier is responsible for predicting one label. The prediction from each classifier in the chain is passed as an additional input to the subsequent classifiers, thereby capturing inter-label relationships effectively. The length of classifier chain gives us the number of labels in the given dataset. Illustrated in Figure 5 are structural design and operational flow of the classifier chain method. Standard multi-label metrics, such as macro accuracy, α - evaluation score and the Hamming loss, which provide metric models of classification performance across multiple label dimensions, are used for model evaluation.

Figure 5-Multi-Label Chain Classifier Diagram



H. Novel Proposed ETRF Feature Engineering Approach

In this context, the proposed ETRF feature extraction method is analyzed from this section, which is a mixture of the trees (ET) and the Random Forests (RF) algorithms. In this study, ETRF is used as a means to extract features that improve the predictive capability of machine learning models for the in identifying genetic disorders. Figure 6 shows the ETRF based feature construction process in terms of overall architecture and data flow. First, the genomic data samples are fed to the ET and RF models separately.

From each model, **class probability predictions** are retrieved and utilized as newly constructed features for subsequent learning models. This fusion of probabilistic outputs from both tree-based ensembles contributes to a more enriched and informative feature set, potentially improving model generalization and classification accuracy.

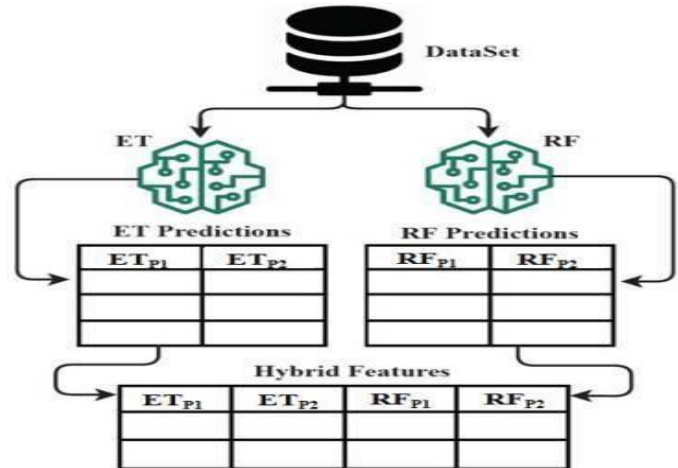


Figure 6 – Class Probability Diagram

V RESULT AND EVALUATION

The genetic disorder prediction model was hosted on Amazon EC2 with anml.c5.xlarge instance, offering a robust and elastic environment appropriate for high through put genomic analysis. The model, using ETRF feature extraction and classifier chaining form multilabel prediction, was benchmarked using standard metrics. High macro accuracy indicated consistent performance across all disorder classes. The F1score indicated a good balance between precision and recall, and the low Hamming loss indicated few mislabel predictions per example. The α evaluating score confirmed the advantage of maintaining label correlations with classifier chaining. In deployment, the EC2 configuration offered consistent and efficient performance with low latency during inference, as well as the ability to handle varying volumes of input data. The overall configuration offered cost effective, accurate, and responsive prediction, making the solution highly appropriate for real world clinical applications involving complex genomic data.

VI EVALUATION METRICS

It is compared based in various important metrics including macro accuracy, α evaluation score, recall, precision, Hamming loss and F1 score. These are all varying aspects of classification and are important while determining the way in which the model is performing with the multi label multi class nature of the dataset. The most important factors utilized while computing these measures are:

- **True Positive:** The Number of Positive Samples Which Are Properly Classified by the Model.
- **True Negative:** Number of correctly predicted negative samples by the model
- **False Positive:** Number of incorrectly predicted negative as positive samples by the model.

Hamming Loss is a way to measure how often a model makes mistakes in multi-label classification. It computes this as incorrect labels out of the total labels and predicts ratio. Essentially, it looks at how many labels were wrongly predicted (Either missing a correct label or assigning an incorrect one) For each instance and then finds the average across all instances.

In simple terms the lower the hamming loss the better the model is at making accurate predictions. A value close to zero means very few mistakes, while a higher value indicates more errors in label predictions

$$\text{HAMMING LOSS} = \frac{1}{N} \sum_{i=1}^N \frac{1}{L} \sum_{j=1}^L |Y_j(i) - \hat{Y}_j(i)|$$

Specifically, the α -evaluation score is a generalization of the Jaccard similarity for the evaluation of the performance of multi-label classification models. On the other hand, it gives more the flexibility as well as a more comprehensive criteria to evaluate how good a learning approach predicts multiple labels for each instance. In simple terms, this score takes into account both false negatives (FN) and false positives (FP) while also considering true positives (TP). The goal is to measure how closely the predicted labels match the actual labels.

$$\alpha\text{-evaluation score} = \left(\frac{\beta M_x + \gamma F_x}{Y_x \vee P_x} \right)^\alpha$$

In Multi-label classification, different evaluation metrics help Measure how well a model is performing.

M_x (false negatives, fn): the number of actual labels That the model failed to predict.

F_x (false positives, fn): the number of incorrect Labels that the model predicted.

Y_x : the total number of true positives (tp) and false

Negatives (fn), representing all actual positive Labels.

P_x : the total number of true positives (tp) and false

Positives (fp), representing all predicted positive Labels.

VII PRECISION AND RECALL

Two widely used measures of assessment in Classification are precision and recall Accuracy is the degree to which the forecast Positive labels are accurate. A high Accuracy means that the model predicts a Positive label, it is usually right. Remember tracks how many of the real Positive labels the model identified correctly. A High recall is where the model performs best at Catching all the labels that apply, even if it Sometimes includes incorrect ones.

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

DEMONEYARATION OF THE COMPARISON OF ACCURACY SCORE WITH APPLIED MACHINE LEARNING METHODS ON A 70:30, 80:20, 85:15, AND 90:10 PROPORTIONATE WAY(DIVISION). AN ANALYSIS OF THE COMPARATIVE ACCURACY RESULT OF THE SUGGESTED TECHNIQUE WITH AND WITHOUT RECOURSE TO IT IN AN IMBALANCED DATASET WITH AND WITHOUT USE OF DATA FIGURE 10 AND FIGURE 11 SHOW AMOUNTS BALANCED. THE ACCURACY RESULTS WITH AND WITHOUT USING THE PROPOSED METHOD WITH DATA BALANCING APPLIED ARE COMPARED.

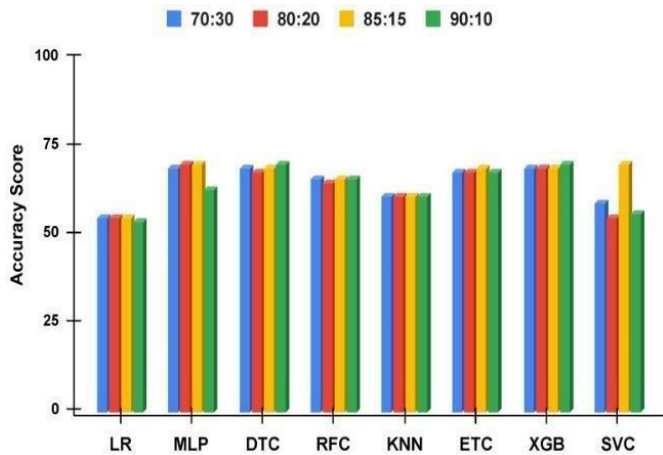


Figure 7 - The graph Employed methods provide a comparative investigation on various data split ratios based on imbalanced datasets without the suggested technique.

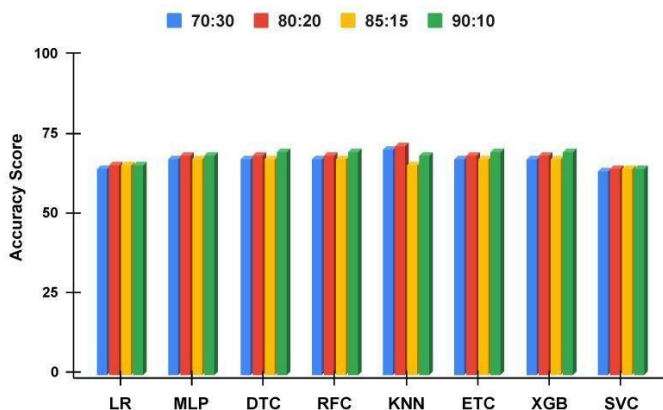


Figure 8 - The graph Utilized Methods Provide a Comparison Of Various data Split Ratios with this Proposed Method When Utilizing imbalanced Data.

VIII DISCUSSION

Utilizing Amazon Web Services (AWS) for genetic disorder prediction has great scalability, performance, and deployment benefits. In this study, AWS EC2 instances were utilized to deploy and execute machine learning models trained on complicated multi-label, multi-class genomic data. The cloud platform allowed for quick model training and prediction, handling large-scale data with low overhead efficiently.

Use of Amazon Sage-Maker and EC2 provided elastic deployment, quick scaling, and high availability. Use of ml.c5xlarge instances was sufficient to provide real-time prediction as well as cost savings. Further, the use of AWS S3 for secure and quick access to data provided smooth data pipelines for prediction and training.

Machine learning pipeline, in particular with the ETRF feature extraction technique and classifier chain strategy, was improved with the computation power and memory of AWS services. Cloud-based distribution facilitated processing of skewed and high-dimensional data with reduced training time. Macro accuracy, precision, recall, and Hamming loss metrics indicated improved results compared to local or constrained environments.

AWS also offers the potential to scale in the future, such as automated retraining on new genomics data, integration with current health data systems, and support for future compliance with health data policies for privacy (e.g., HIPAA). This makes AWS a strong solution for deploying genomics predictive models into research as well as clinic environments.

IX FUTURE SCOPE

Some avenues of work in future are left unsuggested among genetic diseases prediction framework based on machine learning and prediction using AWS. Some of them include employing more sophisticated deep learning models such as CNNs or transformers to possibly intrude more sophisticated patterns of correlation among genetic data and boost predictive accuracy further. Amongst others of utmost priority in area, is clinical integration in real time, which can be scaled up and included in clinics and labs for the immediate prediction while the patient is in consultations.

The combination of lifestyle and environmental variables with genomic information may give rise to a more comprehensive prediction of risk, recognizing the multigenic cause of most genetic disorders. The retraining of models that have been automated with the use of AWS tools such as Sage-maker Pipelines also ensures that the system adapts with the newly arising information and stays relevant and up-to-date with the changing times. Privacy-preserving methodologies such as

federated learning also offer a potential area for extending data across institutions without infringing on sensitive genetic data.

X CONCLUSIONS

This paper describes an AWS scalable infrastructure-based machine learning system for prediction of genetic disorders from genome and genetic data. With careful feature engineering, data balancing, and use of ensemble learning methods—specifically the new ETRF approach—the system makes accurate predictions. Multi-label, multi-class classification was handled using classifier chain management of label dependencies to improve model accuracy.

AWS EC2 deployment via SageMaker offers scalability and real-world relevance, enabling mass-scale and real-time prediction. Experimental results validate the effectiveness of the approach, which shows consistency in the performance metric across several different metrics such as macro accuracy, precision, recall, and F1 score.

In summary, research provides a prospect of integrating cloud infrastructure and machine learning to make genetic research and early diagnosis of inherited diseases a possibility. The system provides an excellent platform to further develop and integrate deep learning, deploy it in real time in clinical setups, and with collaborative learning keeping privacy in mind.

XI REFERENCES

- [1] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32. [Utilized Random Forests, commonly used for the analysis of genetic data.]
- [2] Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794. [Presents XGBoost, used widely in predictive modeling and bioinformatics.]
- [3] Libbrecht, M. W., & Noble, W. S. (2015). Machine learning applications in genetics and genomics. *Nature Reviews Genetics*, 16(6), 321–332. [Describes broad applications of ML in genomics.]
- [4] Min, S., Lee, B., & Yoon, S. (2017). Deep learning in bioinformatics. *Briefings in Bioinformatics*, 18(5), 851–869. [Reviews ML and deep learning approaches in genomics.]
- [5] Eriksson, R., Werling, D. M., et al. (2019). Predicting autism using machine learning on gene expression data. *PLOS ONE*, 14(12), e0226848. [Relevant study on predicting autism based on gene expression.]
- [6] Abid, A., Balin, M. F., & Zou, J. (2019). Concrete Autoencoders for Differentiable Feature Selection and Reconstruction. *International Conference on Machine Learning (ICML)*. [Feature selection methods applicable to genomic data.]
- [7] Amazon Web Services (AWS). (2021). Amazon SageMaker: Developer Guide. <https://docs.aws.amazon.com/sagemaker>
- [8] Shickel, B., Tighe, P. J., et al. (2018). Deep EHR: A survey of recent advances in deep learning techniques for electronic health record (EHR) analysis. *IEEE Journal of Biomedical and Health Informatics*, 22(5), 1589–1604.
- [9] [Healthcare data background on ML models.] Schadt, E. E., et al. (2010). Genetics of gene expression surveyed in maize, mouse, and man. *Nature*, 464(7289), 768–772. [Insights into gene expression patterns and genetic diseases.] Kourou, K., Exarchos, T. P., et al. (2015). Machine learning applications in cancer prognosis and prediction. *Computational and Structural Biotechnology Journal*, 13, 8–17. [Example of ML in predictive genomics for cancer.]