

## **Team – Ansh Pandey & Shajal**

### **Model development**

In the "Anemia Detection with Machine Learning" project, the model development phase was a crucial stage that involved the selection, implementation, and evaluation of various machine learning algorithms. The goal was to develop a model capable of accurately predicting anemia based on clinical features such as gender, hemoglobin levels, and other hematological parameters.

### **Model Selection and Description**

The following machine learning algorithms were considered for the task:

#### **1. Decision Tree (DT):**

- Architecture: Decision Trees are hierarchical models that recursively split the data based on the most significant features. Each internal node represents a decision on a feature, and each leaf node represents an outcome. The model learns to partition the feature space into regions with homogeneous responses.
- Performance: Decision Trees are interpretable and can handle both categorical and continuous data. However, they are prone to overfitting, especially on small datasets, which can lead to poor generalization on unseen data.
- Complexity: The complexity of a Decision Tree is determined by its depth, which directly impacts the model's ability to capture intricate patterns in the data. Shallow trees are less complex but may underfit, while deeper trees are more complex but risk overfitting.
- Computational Requirements: Training a Decision Tree is relatively fast and computationally efficient, even on large datasets. However, pruning techniques may be needed to avoid overfitting.

#### **2. Random Forest (RF):**

- Architecture: Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of their predictions for classification tasks. Each tree is trained on a random subset of the data, with a random subset of features considered at each split.
- Performance: Random Forests are highly effective for classification tasks, providing better generalization than individual Decision Trees. They reduce overfitting by averaging the results of multiple trees, leading to more robust predictions.
- Complexity: The complexity of Random Forests is influenced by the number of trees, their depth, and the number of features considered at each split. While increasing these parameters can improve accuracy, it also increases computational cost.

- Computational Requirements: Training Random Forests can be computationally intensive, particularly with a large number of trees and deep individual trees. However, they can be parallelized, making them scalable to large datasets.

### **3. Logistic Regression (LR):**

- Architecture: Logistic Regression is a linear model that estimates the probability of a binary outcome based on a linear combination of input features. It uses the logistic function to map predicted values to probabilities.
- Performance: Logistic Regression is effective when the relationship between features and the outcome is approximately linear. It is interpretable and can be regularized to prevent overfitting, making it suitable for high-dimensional datasets.
- Complexity: Logistic Regression is less complex compared to other models, with the main complexity arising from the number of features. Regularization techniques like L1 (Lasso) and L2 (Ridge) are used to control complexity and improve generalization.
- Computational Requirements: Logistic Regression is computationally efficient and scales well to large datasets. It requires minimal memory and processing power, making it ideal for quick prototyping and baseline modeling.

### **4. K-Nearest Neighbors (KNN):**

- Architecture: KNN is a non-parametric model that classifies a sample based on the majority class among its k-nearest neighbors in the feature space. It does not involve explicit training; instead, it stores all training samples and uses them for prediction.
- Performance: KNN is intuitive and easy to implement, with performance heavily influenced by the choice of k and the distance metric. It performs well on well-separated data but struggles with noisy or high-dimensional data.
- Complexity: The complexity of KNN arises from the need to compute distances between the test sample and all training samples, which can be computationally expensive as the dataset grows.
- Computational Requirements: KNN requires significant computational resources, particularly for large datasets, as it needs to calculate distances for every prediction. However, optimizations like KD-Trees can reduce computational overhead.

### **5. Support Vector Machine (SVM):**

- Architecture: SVMs are powerful classification models that find the optimal hyperplane separating classes in the feature space. For non-linear data, SVMs use kernel functions to map the data into a higher-dimensional space where a linear separator can be found.

- Performance: SVMs are effective in high-dimensional spaces and are robust to overfitting, particularly when the margin between classes is large. They perform well on both linear and non-linear data, depending on the chosen kernel.
- Complexity: The complexity of an SVM model is influenced by the choice of kernel, the regularization parameter ( $C$ ), and the margin. Non-linear SVMs, particularly with complex kernels, can become computationally intensive.
- Computational Requirements: SVMs can be computationally demanding, especially with large datasets and complex kernels. Training time increases with the number of samples and features, but SVMs often yield high accuracy.

## 6. Gaussian Naive Bayes (NB):

- Architecture: Gaussian Naive Bayes is a probabilistic classifier based on Bayes' theorem, assuming independence between features. It models the distribution of features as Gaussian and calculates the posterior probability for each class.
- Performance: Naive Bayes is simple and efficient, performing well on small datasets or when the independence assumption holds. It is particularly effective for text classification and other high-dimensional datasets.
- Complexity: The complexity of Naive Bayes is low, with the main computational task being the estimation of means and variances for each feature and class. This simplicity makes it less prone to overfitting.
- Computational Requirements: Gaussian Naive Bayes is computationally light, requiring minimal resources for both training and prediction. It is suitable for applications where quick inference is needed.

## Model Training and Evaluation

Each of these algorithms was trained on the preprocessed dataset, and their performance was meticulously evaluated using a variety of metrics:

- Accuracy: Measures the proportion of correctly classified samples.
- Area Under the Curve (AUC): Evaluates the trade-off between the true positive rate and the false positive rate.
- Precision: Assesses the proportion of true positive predictions among all positive predictions.
- Recall: Measures the proportion of true positives among all actual positives.
- F1 Score: The harmonic mean of precision and recall, providing a balanced evaluation of both.
- Kappa Stat: A statistical measure of inter-rater agreement for qualitative (categorical) items, correcting for chance agreement.

These metrics were critical in understanding the strengths and weaknesses of each model, allowing for informed decisions on model selection. Additionally, k-fold cross-validation (with  $k=5$ ) was employed to ensure that the models were not overfitting and that their performance would generalize well to unseen data. This involved repeatedly splitting the dataset into training and validation sets, averaging the performance metrics across all folds to obtain a reliable estimate.

Through this rigorous evaluation process, the best-performing model was identified, setting the stage for the subsequent optimization and tuning phase. The chosen model not only met the project's performance criteria but also demonstrated robustness and reliability, making it suitable for deployment in real-world anemia detection applications.