

A Comprehensive Study of Blood Parameters for Disease Classification

April 15, 2024

Authors: Aryan Rezanezhad(251386495) & Dinithi De Silva(251385811)

Course: Advanced Data Analysis

Instructor: Dr. Camila de Souza

Department: Actuarial Science & Statistics



Contents

Introduction	3
Dataset Description	3
Method	4
Results	4
Exploratory data analysis	4
Analysis of the Response Variable	4
Analysis of the Predictor Variables	5
Statistical Analysis	7
Multinational Logistic Regression Model	7
Multinational Logistic Regression Model with Lasso	7
LDA	8
QDA	8
Classification Tree	9
Random Forest	10
K-Nearest Neighbors (KNN)	10
Naïve Bayes Model	10
Conclusion	11
Appendix	12
Tables	12
Rcodes	14
R Code Table of Contents	14
References	22

Introduction

In health care and medical diagnostics, the analysis of blood test results plays a crucial role in the early detection of disease. At present, advancements in technology have significantly enhanced the capabilities of medical diagnostics. These technological improvements include more sophisticated laboratory equipment as well as and software that can analyze complex data sets more efficiently. This means that diseases can be identified more quickly and accurately, allowing for timely intervention and better management of patients' health outcomes.

Doe and Smith (2020) highlighted the importance of using results of blood tests to help doctors detect diseases early. They suggested that combining these data with computer methods could improve these predictions. In 2021, Smith discusses how doctors utilize predictions from computational tools in diagnosing diseases and he emphasizes that while these tools are beneficial, it's crucial for doctors to also rely on their professional judgment, recognizing the unique aspects of each patient's case. These studies collectively indicate that data and computer technology are increasingly integral to healthcare. They enhance doctors' ability to understand patients' health more thoroughly and enable the creation of personalized treatment plans.

This project employs advanced statistical methods to classify the certain diseases based on several blood parameters and compare the performances of each methods. This could lead to more precise and early diagnosis, significantly enhancing disease management and patient care outcomes.

Dataset Description

The dataset, sourced from Kaggle (Aboelnaga, E. (2022)) consists of 24 blood parameters along with a categorical variable (Disease) indicating four disease types (Diabetes, Thalassemia, Anemia, Thrombocytopenia) or a healthy status. The dataset contains 150 entries of pre-processed data, and the predictors within it have been scaled to a range between 0 and 1. This study employs a supervised learning approach, where blood parameters are utilized as predictor variables, and the presence of disease is designated as the response variable for the analysis. Detailed description of variables listed in the table 3 (see Appendix).

Method

The study comprises two primary sections: exploratory analysis and statistical analysis. We initiated the exploratory analysis by analyzing the response variable. Then We analyzed explanatory variables through histograms and statistical summaries to see the distribution, central tendencies, and variability of the data. Then, we obtained statistical summaries of all the predictor variables against the response to understand how the distributions vary across the response categories. Next, the relationships between predictor variables have been assessed using a correlation matrix (Pearson’s correlation coefficient) to identified the potential multicollinearity issues among the predictor variables.

The statistical analysis started by dividing the dataset into training (80%) and test (20%) sets. Our first attempt was the multinomial logistic regression model with and withoutlasso penalization and compare the performance. Next, attempt was the linear discriminant analysis (LDA) and Quadratic discriminant analysis (QDA) but due to the violation of the required assumptions it failed to implement both models.Next, we implemented classification tree and random forest and compare the performances. Finally we employed both the k-nearest neighbor (KNN) method and the Naïve Bayes Model to assess and compare their respective performances

Results

Exploratory data analysis

Analysis of the Response Variable

In Table 1, the frequency distribution and proportions across disease categories are depicted. It is noteworthy that the Healthy group slightly marginally outnumbers the others, while Thromboc exhibits the lowest frequency.When comparing proportions, no single disease category dominates the sample, ensuring an equitable distribution that provides a solid basis for unbiased predictive modeling.

Table 1: Frequency distribution of Disease variable

Disease	Frequency	Proportion
Anemia	30	0.20
Diabetes	34	0.23
Healthy	35	0.23
Thalasse	27	0.18
Thromboc	24	0.16

Analysis of the Predictor Variables

a) Univariate Analysis The dataset comprises several key biochemical and hematological parameters, each potentially indicative of underlying health status. The distribution of these variables provides an initial insight into the nature of the data. (Figure 1). Notably, the histogram for glucose levels exhibits a right-skewed distribution, suggesting that patients with higher glucose levels may have diabetic conditions. The Red Blood Cell count demonstrates a bimodal distribution, indicating the potential existence of distinct subgroups within our sample, possibly differentiating between anemic and non-anemic individuals. Correspondingly, similar bimodal patterns observed in Hematocrit and Hemoglobin distributions support this observation and justify a more detailed investigation into their association with conditions related to anemia.

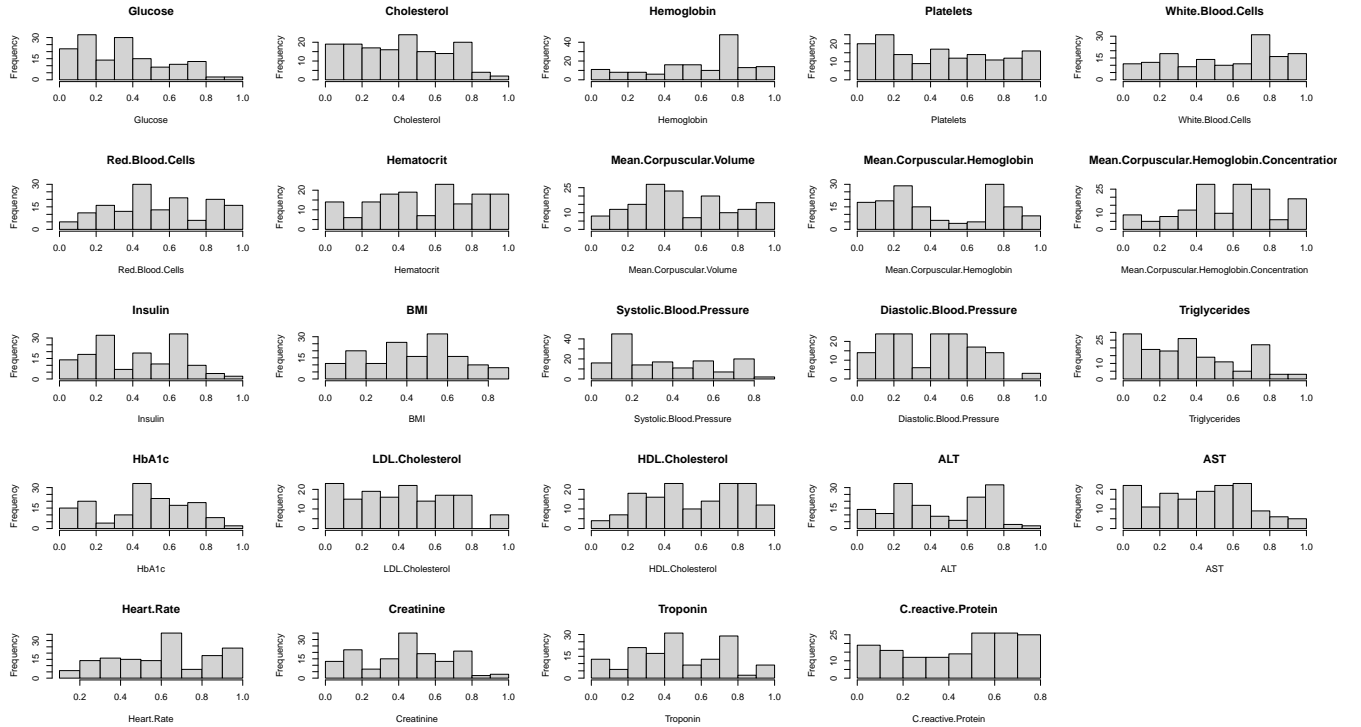


Figure 1: Histograms for the predictors

Statistical summaries of the predictors (see Appendix, Table 4) provides a detailed statistical overview of various blood parameters within our study population. Glucose levels average at 0.349 with moderate variability, and its slightly right-skewed distribution has a median of 0.353. Cholesterol and hemoglobin illustrate different distribution characteristics; cholesterol levels slightly exceed the mean with a median of 0.412, whereas hemoglobin displays a higher median than the mean. Additionally, white and red blood cells, along with hematocrit, present higher baseline values, generally exceeding 0.5, with red blood cells and hematocrit showing lower variability.

Next, We obtained the standard deviations across the disease group (see Appendix, Tables 5, 6, 7). We can observed that Mean Corpuscular Volume and Mean Corpuscular Hemoglobin display higher variability within the Thalassemia group. This finding is consistent with the known hematological manifestation of Thalassemia, which affects the size and hemoglobin content of red blood cells. In contrast, variables such as C-reactive Protein exhibit lower variability, suggesting a more stable measurement unaffected by the specific disease state.

b) Bivariate Analysis Figure 2 presents the correlation matrices for each combination of predictor variables. It shows that there is somewhat strong positive correlations between the glucose and the BMI level, which suggests potential impact of body weight on glucose metabolism, highlighting the importance of weight management in the prevention and control of diabetes. Furthermore, This visual correlation plot confirmed that there is no potential multicollinearity issues, which is crucial for model selection and ensuring the independence of predictors in regression analyses.



Figure 2: Correlation pair plot

Statistical Analysis

Under the statistical analysis, we utilized various classification techniques to classify the diseases based on blood parameters. Each model was trained using the same training dataset and subsequently evaluated using a distinct test set. Our primary focus was on assessing the misclassification error rate to compare the performance of each model.

Multinational Logistic Regression Model

Our initial approach employed a multinomial logistic regression model, designating 'Healthy' as the reference category. This model enables direct comparisons of disease likelihoods relative to a baseline of health, facilitating a clear understanding of how disease probabilities differ from the healthy state.

The model's coefficients reveal the relative influence of each predictor. For instance, Glucose and Cholesterol exhibit considerable effects on the likelihood of being classified as Diabetic versus Healthy. The magnitude of these coefficients suggests a strong relationship between these predictors and disease status, aligning with medical literature that implicates these factors in metabolic syndromes.

Interestingly, the variables related to blood cell characteristics, such as Red Blood Cells and Hematocrit, show substantial weights in the model, possibly underlining their diagnostic value in distinguishing among disease states. These findings reflected through the coefficients corresponding to Anemia and Thalassemia, diseases intrinsically related to blood cell properties. (Figure 3)

The error rate for the multinational logistic model is approximately 14%, indicating a somewhat high degree of accuracy in distinguishing between diseases based on blood test data.

```
## [1] "Misclassification Error Rate: 0.14"
```

Multinational Logistic Regression Model with Lasso

Next, We implemented a lasso regression model, which introduces a penalty on the absolute size of the coefficients. This regularization technique promotes a sparse model by shrinking coefficients of less influential predictors towards zero, thereby enhancing interpretability.

The model fitting procedure employs cross-validation to select the penalty strength, ensuring that the chosen model generalizes well to unseen data. The results from this model underscore the predictive power of some markers over others, with variables like Mean Corpuscular Hemoglobin and BMI surfacing as significant after the lasso penalty is applied.

In evaluating model efficacy, we consider the misclassification error rate, a straightforward metric quantifying the proportion of incorrect predictions. The lasso regression model exhibits a misclassification error rate of 10%, denoting an approximately 90% accuracy in classification on the test set, a testament to the model's precision in disease state prediction.

```
## [1] "Misclassification Error Rate: 0.1"
```

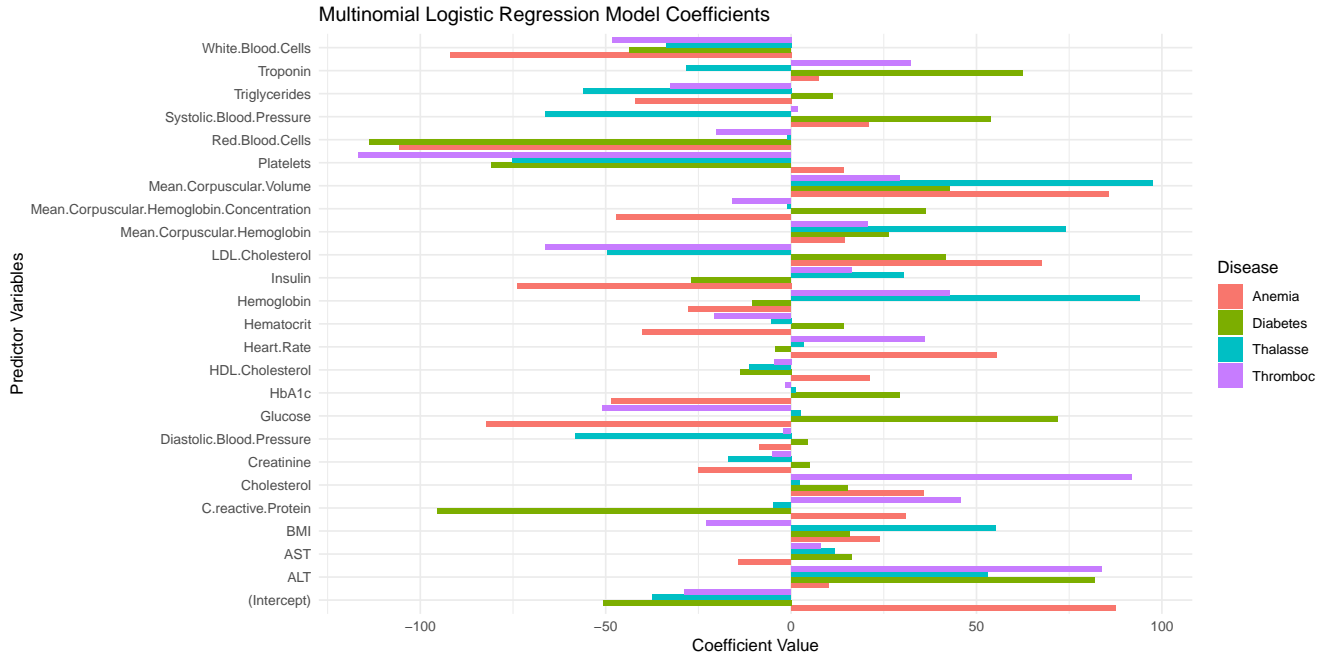


Figure 3: Multinomial Logistic Regression Model Coefficients

LDA

Linear Discriminant Analysis (LDA) is a statistical method used for classification and dimensionality reduction. LDA assumes that the predictor variables are normally distributed within each class and have the same variance-covariance matrix across the classes (homoscedasticity). These assumptions are critical for the accurate performance of the model.

During the EDA, we analyzed the standard deviations for each predictor variable segmented by disease classification. It was complex in understanding how the variances fluctuated across different disease conditions. The results, as we discussed in the EDA part of our report, indicated noticeable differences in variabilities among groups, suggesting a potential violation of the homoscedasticity assumption required by LDA.

Alternative approach is Quadratic Discriminant Analysis (QDA), which does not assume equal covariance matrices, may be more appropriate.

QDA

In our efforts to employ advanced statistical methods to understand and predict disease states from our dataset, we considered Quadratic Discriminant Analysis (QDA) as a potential model. QDA is favored for its ability to model and handle datasets where classes have distinct covariance structures. However, during the implementation phase, we encountered with a issue that prevented the application of QDA to our current dataset.

Our evaluation of Quadratic Discriminant Analysis (QDA) for classifying disease states within our dataset revealed a significant impediment due to high dimensionality. The application of QDA was hindered by a rank deficiency in the covariance matrix of the “Healthy” group, where the number of predictor variables exceeds the number of available observations. This mathematical

constraint prevents the inversion of the covariance matrix, a critical step in QDA, rendering the method unsuitable under our current dataset configuration.

Classification Tree

In our study, a classification tree was constructed to predict various disease states based on a set of clinical measurements. Utilizing the rpart package in R, we crafted a model that simplifies the complexity of medical diagnoses into an interpretable series of decisions. The decision tree, depicted in a detailed diagram, outlines the decision-making process by utilizing key factors such as mean corpuscular volume, Hemoglobin levels, and other blood parameters to categorize patients into groups including Healthy, Anemia, Thromboc, and Diabetes. (Figure 4)

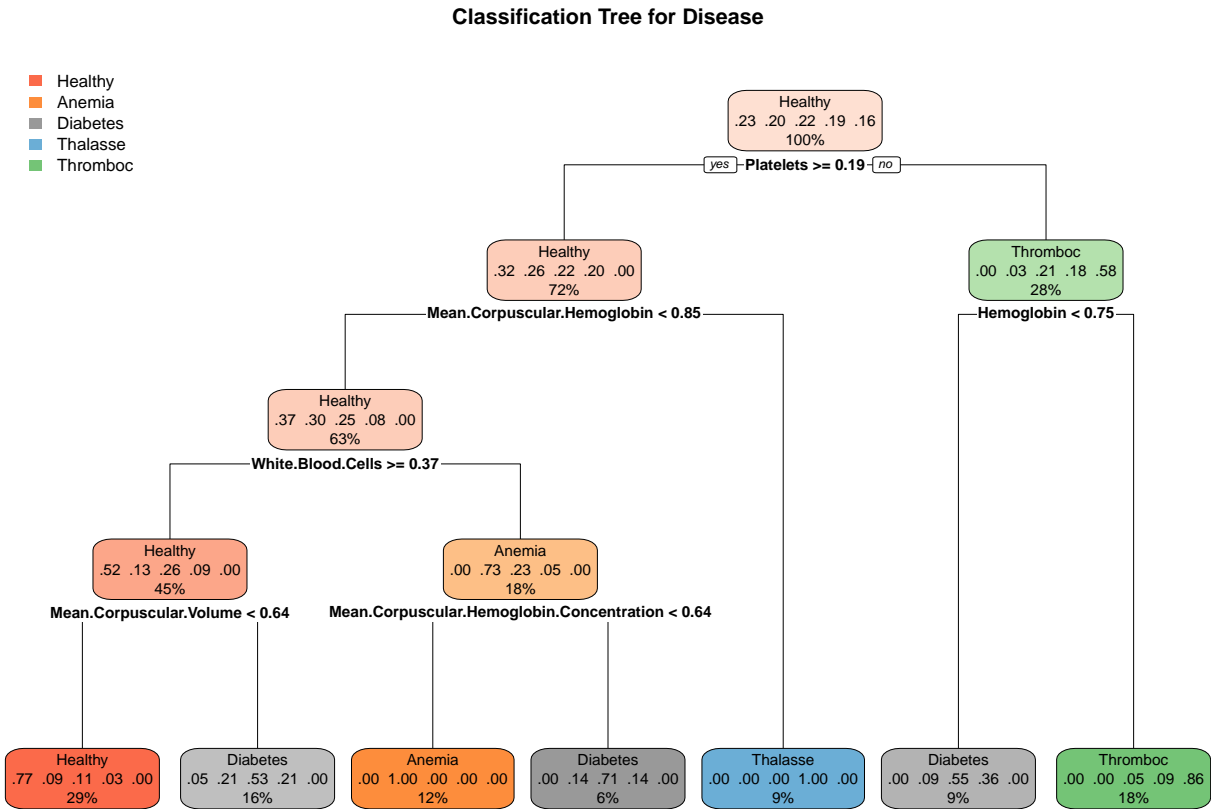


Figure 4: Classification Tree for Disease

The classification tree model indicated a misclassification error rate of 27%, as computed on the test dataset. This error rate represents the proportion of patients whose disease status was incorrectly predicted by the model. While the classification tree provides an intuitive understanding of the data through its clear decision rules, the error rate also underscores the limitations in predictive performance. Such findings prompt consideration of model refinement through techniques like pruning or exploring ensemble methods that may yield more accurate predictions while maintaining the interpretability that is vital in clinical settings.

```
## [1] "Misclassification Error Rate: 0.27"
```

Random Forest

The random forest model demonstrated a misclassification error rate of 3%, as evaluated on the test dataset. This rate quantifies the percentage of patients whose disease status was inaccurately predicted by the model. Although the random forest offers robustness through its aggregation of decision trees, thus typically improving prediction accuracy over a single decision tree, this error rate highlights its precision in identifying disease states. The relatively low error rate suggests effective disease classification, reinforcing the model's utility in clinical applications.

```
## [1] "Misclassification Error Rate: 0.03"
```

K-Nearest Neighbors (KNN)

We used the K-Nearest Neighbors (KNN) method with cross-validation (CV) value to predict the correct disease from medical data. KNN looks at what's known about a patient and finds other patients who had similar test results. To make sure to every test result was given equal importance, we first made all the numbers more comparable to each other. We also used a special process to pick the best number of similar patients (neighbors) to look at for making our predictions.

After running the KNN model, it achieved a correctness rate of approximately 67%, indicating an error rate of approximately 33%. This gives us a good starting point, showing us that the method can indeed spot patterns and give us useful predictions. Nevertheless, there may have more enhancement for this model. Consideration may be given to altering the methodology employed for determining the number of neighbors other than CV method or refining the dataset to potentially optimize the model's performance in subsequent iterations

```
## [1] "Misclassification Error Rate: 0.33"
```

Naïve Bayes Model

In our data analysis project, we explored the use of the Naïve Bayes classifier, a straightforward yet effective machine learning technique, to predict various diseases based on a range of clinical indicators. After fitting the model to a portion of our data reserved for training, we then tested its predictive power on a separate test set. The results were promising, with the Naïve Bayes model achieving a 33% misclassification error rate, suggesting that it correctly predicted the disease status 67% of the time.

The strength of the Naïve Bayes classifier lies in its simplicity and assumption that the features in the dataset are independent of each other given the disease status. Despite this simplification, which may not hold in real-world settings where features can be interrelated, the model performed well. This performance is especially notable given the complex nature of medical diagnosis, where even a small improvement in predictive accuracy can have significant implications. However, for more nuanced interpretations and potentially improved results, further fine-tuning of the model or the use of more sophisticated algorithms might be considered for future analyses.

```
## [1] "Misclassification Error Rate: 0.33"
```

Conclusion

In our analysis of various classification models for disease classification, several key findings emerged. The multinomial logistic regression model, when utilizing LASSO penalty, demonstrated a notably lower misclassification error rate of approximately 10%, compared to approximately 14% when fitted without LASSO. This result underscores its high accuracy in distinguishing between diseases based on blood test data. Notably, variables related to blood cell characteristics, such as Red Blood Cells and Hematocrit, exhibited substantial weights in the model, suggesting their diagnostic significance in identifying disease states associated with conditions like Anemia and Thalassemia.

However, Linear Discriminant Analysis (LDA) posed challenges due to assumptions of homoscedasticity, as indicated by notable differences in variabilities among disease groups during exploratory data analysis. Quadratic Discriminant Analysis (QDA) was considered as an alternative but encountered issues with rank deficiency in covariance matrices, rendering it unsuitable for our dataset.

The classification tree model exhibited a misclassification error rate of 27%, highlighting its interpretability yet indicating limitations in predictive performance. Conversely, the random forest model demonstrated a lower error rate of 3%, which highlights its accuracy in classifying the disease states

Upon running the KNN model, an initial correctness rate of approximately 67% was achieved, suggesting potential for pattern recognition. However, further enhancements are warranted, possibly through alternative methodologies for determining the number of neighbors or dataset refinement.

Finally, the Naïve Bayes classifier showed promise with a 33% misclassification error rate, indicating correct prediction of disease status 67% of the time.

In conclusion, each model exhibited distinct strengths and limitations in disease classification. The choice of model should consider factors such as interpretability, predictive accuracy, and computational efficiency, with ongoing refinement through methodological adjustments and dataset optimization to enhance performance in clinical decision-making scenarios. In terms of the prediction accuracy, the random forest model demonstrates superior accuracy with a low error rate of 3%, making it the most effective model for predicting disease states from the provided dataset. (Table 2)

Table 2: Misclassification Error Rates of Different Classification Models	
Model	Misclassification Error Rate
Multinomial Logistic Regression (with LASSO)	10%
Multinomial Logistic Regression (without LASSO)	14%
Linear Discriminant Analysis (LDA)	Not Applicable
Quadratic Discriminant Analysis (QDA)	Not Applicable
Classification Tree	27%
Random Forest	3%
K-Nearest Neighbors (KNN)	33%
Naïve Bayes	33%

Appendix

Tables

Table 3: Dataset Variables

Variable Description	Variable Type
Glucose (mg/dL)	Continuous
Cholesterol (mg/dL)	Continuous
Hemoglobin (g/dL)	Continuous
Platelets (per microliter blood)	Continuous
White Blood Cells (per cubic millimeter of blood)	Continuous
Red Blood Cells (million cells per microliter of blood)	Continuous
Hematocrit (percentage of blood volume)	Continuous
Mean Corpuscular Volume (average volume of red blood cells)	Continuous
Mean Corpuscular Hemoglobin (average amount of hemoglobin in a red blood cell)	Continuous
Mean Corpuscular Hemoglobin Concentration:(grams per deciliter)	Continuous
Insulin (microU/mL)	Continuous
BMI (kg/m ²): Body Mass Index	Continuous
Systolic Blood Pressure (mmHg)	Continuous
Diastolic Blood Pressure (mmHg)	Continuous
Triglycerides (mg/dL)	Continuous
HbA1c (percentage)	Continuous
LDL Cholesterol (mg/dL)	Continuous
HDL Cholesterol (mg/dL)	Continuous
ALT: Alanine Aminotransferase (Unit per liter)	Continuous
AST (U/L): Aspartate Aminotransferase (Unit per liter)	Continuous
Heart Rate (beats per minute)	Continuous
Creatinine (mg/dL)	Continuous
Troponin (ng/mL)	Continuous
C-reactive Protein (mg/L)	Continuous
Disease (Healthy, Diabetes, Thalassemia, Anemia, Thrombocytopenia)	Categorical

Table 4: Summary of the Predictors

Variable	mean	SD	median	IQR	min	max
Glucose	0.3486714	0.2430851	0.3527283	0.3608228	0.0109939	0.9684602
Cholesterol	0.4075077	0.2427988	0.4123852	0.3990268	0.0121394	0.9050264
Hemoglobin	0.5963739	0.2681275	0.6818528	0.3105240	0.0030213	0.9833058
Platelets	0.4563661	0.3066210	0.4159568	0.5518270	0.0125938	0.9993931
White.Blood.Cells	0.5482071	0.2880032	0.6186836	0.5335831	0.0101385	0.9907857
Red.Blood.Cells	0.5350428	0.2593967	0.5239267	0.3643961	0.0445651	1.0000000
Hematocrit	0.5466259	0.2796318	0.5428600	0.4501717	0.0117717	0.9775203
Mean.Corpuseular.Volume	0.5073132	0.2632431	0.4607011	0.3889384	0.0469415	0.9952629
Mean.Corpuseular.Hemoglobin	0.4507419	0.3019482	0.3143288	0.5531111	0.0005536	0.9632348
Mean.Corpuseular.Hemoglobin.Concentration	0.5721316	0.2557413	0.6272867	0.3077711	0.0069470	0.9755857
Insulin	0.4195688	0.2420510	0.4155118	0.4341323	0.0341291	0.9667836
BMI	0.4302201	0.2336386	0.4039203	0.3360820	0.0145956	0.8982099
Systolic.Blood.Pressure	0.3519395	0.2342432	0.2990437	0.4029925	0.0059883	0.8290996
Diastolic.Blood.Pressure	0.4014083	0.2345565	0.4243051	0.4021625	0.0055786	0.9346175
Triglycerides	0.3834101	0.2555339	0.3886270	0.4361552	0.0106123	0.9736794
HbA1c	0.4631550	0.2511161	0.4895137	0.4352284	0.0162556	0.9502181
LDL.Cholesterol	0.4040407	0.2573274	0.4088922	0.4048151	0.0330373	0.9838265
HDL.Cholesterol	0.5603782	0.2553466	0.5626217	0.4209567	0.0395053	0.9772271
ALT	0.4525754	0.2607057	0.3922509	0.4328867	0.0071863	0.9425488
AST	0.4297123	0.2480985	0.4552624	0.3775721	0.0130125	0.9944600
Heart.Rate	0.6060094	0.2499545	0.6121878	0.4213373	0.1145503	0.9968728
Creatinine	0.4377443	0.2342008	0.4252962	0.4210546	0.0763325	0.9259243
Troponin	0.4808100	0.2463520	0.4683440	0.4135965	0.0074902	0.9728028
C.reactive.Protein	0.4437758	0.2362612	0.5182443	0.4037515	0.0048673	0.7979059

Table 5: SD of Each Variable Across Disease Groups - Part 1

Disease	Glucose	Cholesterol	Hemoglobin	Platelets	White.Blood.Cells	Red.Blood.Cells
Healthy	0.2305791	0.2204935	0.1957583	0.1941311	0.1841050	0.2347559
Anemia	0.1862419	0.2036432	0.2679817	0.3117764	0.2507779	0.2984934
Diabetes	0.2965607	0.3011817	0.3355210	0.3301893	0.2813515	0.2598796
Thalasse	0.2193233	0.1959411	0.2668046	0.2771660	0.2296850	0.2219584
Thromboc	0.1258863	0.1084362	0.0153098	0.0689885	0.3088241	0.1676606

Table 6: SD of Each Variable Across Disease Groups - Part 2

Disease	Hematocrit	Mean.Corpuseular.Volume	Mean.Corpuseular.Hemoglobin	Mean.Corpuseular.Hemoglobin.Concentration	Insulin	BMI	Systolic.Blood.Pressure
Healthy	0.2697968	0.1762608	0.3057492	0.1672030	0.2513395	0.1857696	0.1989325
Anemia	0.3138148	0.3033710	0.2378729	0.2198056	0.2271171	0.2342215	0.2485292
Diabetes	0.2945140	0.2742189	0.3100812	0.2783050	0.2364421	0.2720255	0.2381060
Thalasse	0.2257225	0.3085976	0.3619884	0.3486779	0.2196849	0.2496707	0.2209132
Thromboc	0.1652991	0.0368136	0.2261115	0.1190782	0.2324222	0.1184450	0.2788758

Table 7: SD of Each Variable Across Disease Groups - Part 3

Disease	Diastolic.Blood.Pressure	Triglycerides	HbA1c	LDL.Cholesterol	HDL.Cholesterol	ALT	AST	Heart.Rate	Creatinine
Healthy	0.2455835	0.2571616	0.2397051	0.1729278	0.1775510	0.2242089	0.2144057	0.2558157	0.2332777
Anemia	0.2122082	0.2118387	0.2399524	0.2342866	0.1862880	0.1855852	0.2242483	0.2558081	0.1572440
Diabetes	0.2539544	0.3333171	0.2799410	0.3193529	0.2988492	0.2553456	0.3407486	0.2841400	0.2973660
Thalasse	0.2552282	0.1822867	0.2762595	0.1646382	0.2893404	0.2736918	0.2237212	0.2152382	0.2213999
Thromboc	0.1476657	0.1573757	0.1015483	0.3294890	0.1042897	0.2342772	0.1904076	0.1614171	0.1576377

Rcodes

R Code Table of Contents

• Libraries	15
• Read the data	15
• Frequency table for disease variable	15
• Descriptive Analysis - predictor variables	16
• Correlation pair plot	16
• Multinomial Logistic Regression Model	16
• Multinomial Logistic Regression Model Coefficients	17
• Multinomial Logistic Regression Model with Lasso	17
• Classification Tree for Disease Prediction	18
• Random Forest	19
• K-Nearest Neighbors (KNN)	19
• Naïve Bayes model	19
• Summary of the Predictors	20
• SD of each variable across the groups in response	20

```

1 library(tidyverse)
2 library(ISLR)
3 library(reshape2)
4 library(ggplot2)
5 library(faraway)
6 library(rsm)
7 library(dplyr)
8 library(purrr)
9 library(knitr)
10 library(kableExtra)
11 library(nnet)
12 library(glmnet)
13 library(tidyr)
14 library(caret)
15 library(randomForest)
16 library(rpart)
17 library(rpart.plot)
18 library(class)
19 library(e1071)

```

Listing 1: Libraries

```

1 # Read the data
2 blood <- read.csv("bloodEx3.csv")
3 str(blood)
4
5 # Divide the training and test set
6 set.seed(4)
7 n <- nrow(blood)
8 ind <- sample(nrow(blood), 0.8*n)
9 train.set <- blood[ind,]
10 test.set <- blood[-ind,]

```

Listing 2: Read the data

```

1 #| label: tbl-1
2 #| tbl-cap: "Frequency table for disease variable "
3 #| tbl-height: 8
4 #| tbl-width: 12
5
6 # Frequency table for disease variable
7 frequency_table <- table(blood$Disease)
8 frequency_df <- as.data.frame(frequency_table)
9 names(frequency_df) <- c("Disease", "Frequency")
10
11 # Proportion table
12 prop_table <- prop.table(frequency_table)
13 prop_df <- as.data.frame(prop_table)
14 names(prop_df) <- c("Disease", "Proportion")
15 prop_df$Proportion <- round(prop_df$Proportion, 2)
16
17 # Merge frequency and proportion data frames
18 combined_df <- merge(frequency_df, prop_df, by = "Disease")
19 combined_df <- as.data.frame(combined_df)
20
21 # Create frequency table with kable
22 kable(combined_df, caption = "Frequency distribution of Disease variable",

```

```

23     booktabs = TRUE) |>
24     kable_styling(bootstrap_options = c("striped", "hover"))

```

Listing 3: R code for the Table 1

```

1  #| label: fig-1
2  #| fig-cap: "Histograms for the predictors"
3  #| fig-height: 9
4  #| fig-width: 16
5
6  # Descriptive Analysis - predictor variables
7  par(mfrow = c(5,5))
8  for (variable_name in names(blood[-25])) {
9    hist(blood[[variable_name]], main = variable_name, xlab = variable_name)}

```

Listing 4: R code for the Figure 1

```

1  #| label: fig-2
2  #| fig-cap: "Correlation pair plot"
3  #| fig-height: 14
4  #| fig-width: 16
5
6  # Calculate the correlation matrix for the first 24 variables
7  cor_matrix <- cor(blood[, 1:24], use = "complete.obs")
8
9  # Convert the correlation matrix into a tidy data frame
10 cor_data <- as_tibble(cor_matrix, rownames = "Variable1") |> # Use as_tibble for
    rownames
11 pivot_longer(cols = -Variable1, names_to = "Variable2",
12               values_to = "Correlation") |>
13               # Optionally, filter out self-correlations
14               filter(Variable1 != Variable2) |>
15               mutate(AbsoluteCorrelation = abs(Correlation),
16                      Variable2 = stringr::str_wrap(Variable2, width = 15)) |>
17               arrange(desc(AbsoluteCorrelation))
18
19 # Create a correlation heatmap
20 ggplot(cor_data, aes(Variable2, Variable1, fill = Correlation)) +
21   geom_tile(height = 0.9) +
22   scale_fill_gradient2(low = "blue", mid = "white", high = "red",
23                       limits = c(-1, 1)) +
24   theme_minimal() +
25   theme(axis.text.x = element_text(angle = 90, vjust = 0.5,
26                                     hjust = 1, size = 6),
27         axis.text.y = element_text(size = 6)) +
28   coord_flip()

```

Listing 5: R code for the Figure 2

```

1  # Multinational Logistic Regression Model
2  # Convert the response variable to a factor
3  train.set$Disease <- as.factor(train.set$Disease)
4
5  # Prepare the data for glmnet
6  x <- model.matrix(Disease ~ . -1, data = train.set) # Predictor matrix, removing
    the intercept
7  y <- train.set$Disease
8

```



```

9 # Fit a multinomial logistic regression model without LASSO penalty
10 set.seed(2)
11 model <- glmnet(x, y, family = "multinomial")
12
13 # Prepare the predictor matrix for test_data similar to how it was done for the
  training data
14 x_test <- model.matrix(Disease ~ . -1, data = test.set)
15
16 # Predicting on the test data
17 predictions <- predict(model, newx = x_test, type = "class")
18
19 # Calculate the misclassification error rate
20 actual <- test.set$Disease # Actual response variable from the test set
21 misclass_error_rate <- mean(predictions != actual)
22
23 # Print the misclassification error rate
24 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
  ))

```

Listing 6: R code for the Multinomial Logistic Regression Model

```

1 #| label: fig-3
2 #| fig-cap: "Multinomial Logistic Regression Model Coefficients"
3 #| fig-height: 6
4 #| fig-width: 12
5 blood$Disease <- as.factor(blood$Disease)
6 set.seed(2)
7 # Fit the model
8 multinom_model <- multinom(Disease ~ ., data = blood)
9
10 # Extract coefficients
11 model_coef <- coef(multinom_model)
12
13 # Transform the coefficients to a tidy data frame for plotting
14 coef_df <- as.data.frame(t(model_coef))
15 coef_df$Predictor <- rownames(coef_df)
16 coef_df <- coef_df |>
17   gather(key = "Disease", value = "Coefficient", -Predictor)
18
19 # Plot the coefficients using ggplot2
20 ggplot(coef_df, aes(x = Predictor, y = Coefficient, fill = Disease)) +
21   geom_bar(stat = "identity", position = position_dodge()) +
22   theme_minimal() +
23   labs(title = "Multinomial Logistic Regression Model Coefficients",
24        x = "Predictor Variables",
25        y = "Coefficient Value") +
26   coord_flip() # Flipping coordinates for a horizontal bar plot

```

Listing 7: R code for the Figure 3

```

1 # Multinomial Logistic Regression Model with Lasso
2 # Convert the response variable to a factor
3 train.set$Disease <- as.factor(train.set$Disease)
4 set.seed(2)
5 # Prepare the data for glmnet
6 x <- model.matrix(Disease ~ . -1, data = train.set) # Predictor matrix, removing
  the intercept

```

```

7 y <- train.set$Disease
8
9 # Convert y to a matrix
10 y <- model.matrix(~y - 1)
11
12 # Fit a multinomial logistic regression model with LASSO penalty
13 set.seed(2)
14 cv_glmnet <- cv.glmnet(x, y, family = "multinomial", alpha = 1, parallel = TRUE)
15
16 # Fit the final model using the best lambda
17 final_model <- glmnet(x, y, family = "multinomial", alpha = 1)
18
19 # To see the coefficients
20 #coef(final_model)
21
22 # Prepare the predictor matrix for test_data similar to how it was done for the
    training data
23 x_test <- model.matrix(Disease ~ . -1, data = test.set)
24
25 # Predicting on the test data
26 predictions <- predict(final_model, newx = x_test, type = "class",
27                        s = cv_glmnet$lambda.min)
28 predictions_char <- as.character(predictions)
29
30 # Remove the 'y' prefix from predictions
31 predictions_char <- sub("~y", "", predictions_char)
32
33 # Now convert back to factor, using the original levels of the Disease variable
34 predictions_factor <- factor(predictions_char)
35
36 # Calculate the misclassification error rate
37 actual <- test.set$Disease # Actual response variable from the test set
38 misclass_error_rate <- mean(predictions_factor != actual)
39
40 # Print the misclassification error rate
41 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
    ))

```

Listing 8: R code for the Multinational Logistic Regression Model with Lasso

```

1 #| label: fig-4
2 #| fig-cap: "Classification Tree for Disease"
3 #| fig-height: 8
4 #| fig-width: 12
5
6 n <- nrow(blood)
7 set.seed(2)
8 ind <- sample(nrow(blood), 0.8*n)
9 train <- blood[ind,]
10 test <- blood[-ind,]
11 # Fit the classification tree
12 tree_model <- rpart(Disease ~ ., data = train, method = "class")
13 # Plotting the tree
14 rpart.plot(tree_model, main="Classification Tree for Disease Prediction")
15
16 # Predict on the test set
17 predictions <- predict(tree_model, newdata = test, type = "class")

```

```

18
19 # Calculate the misclassification error rate
20 actual <- test$Disease
21 misclass_error_rate <- mean(predictions != actual)
22
23 # Print the misclassification error rate
24 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
  ))

```

Listing 9: R code for the Figure 4 and Classification Tree

```

1 # Random Forest
2 rf_model <- randomForest(Disease ~ ., data = train.set, ntree = 500)
3 test_predictions <- predict(rf_model, newdata = test.set)
4 misclass_error_rate <- mean(test_predictions != test.set$Disease)
5 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
  ))

```

Listing 10: R code for the Random Forest Model

```

1 # KNN
2 n <- nrow(blood)
3 set.seed(2)
4 ind <- sample(nrow(blood), 0.8*n)
5 train <- blood[ind,]
6 test <- blood[-ind,]
7
8 # Normalize the data
9 normalize <- function(x) {return ((x - min(x)) / (max(x) - min(x)))}
10
11 train_norm <- as.data.frame(lapply(train[, -ncol(train)], normalize))
12 test_norm <- as.data.frame(lapply(test[, -ncol(test)], normalize))
13
14 # Include the Disease column back
15 train_norm$Disease <- train$Disease
16 test_norm$Disease <- test$Disease
17
18 # Determine the best value of k using cross-validation
19 set.seed(123)
20 ctrl <- trainControl(method = "cv", number = 10)
21 knn_fit <- train(Disease ~ ., data = train_norm, method = "knn",
22               trControl = ctrl, preProcess = c("center", "scale"))
23 best_k <- knn_fit$bestTune$k
24
25 # Fit the KNN model
26 knn_pred <- knn(train = train_norm[, -ncol(train_norm)],
27               test = test_norm[, -ncol(test_norm)],
28               cl = train_norm$Disease, k = best_k)
29
30 # Calculate the misclassification error rate
31 actual <- test_norm$Disease
32 misclass_error_rate <- mean(knn_pred != actual)
33
34 # Print the misclassification error rate
35 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
  ))

```

Listing 11: R code for the K-Nearest Neighbors Model

```

1 # Fit the Naïve Bayes model
2 nb_model <- naiveBayes(Disease ~ ., data = train)
3 # Predict on the test set
4 predictions <- predict(nb_model, newdata = test)
5
6 # Calculate the misclassification error rate
7 actual <- test$Disease
8 misclass_error_rate <- mean(predictions != actual)
9
10 # Print the misclassification error rate
11 print(paste("Misclassification Error Rate:", round(misclass_error_rate, digits = 2)
  ))

```

Listing 12: R code for the Naïve Bayes model

```

1 #| label: tbl-3
2 #| tbl-cap: "Summary of the Predictors"
3 #| tbl-height: 4
4 #| tbl-width: 14
5
6 # Custom summary function for each column
7 custom_summary <- function(x, name) {
8   data.frame(
9     Variable = name,
10    mean = mean(x, na.rm = TRUE),
11    SD = sd(x, na.rm = TRUE),
12    median = median(x, na.rm = TRUE),
13    IQR = IQR(x, na.rm = TRUE),
14    min = min(x, na.rm = TRUE),
15    max = max(x, na.rm = TRUE)
16  )
17 }
18 summary_df <- map_dfr(1:24, function(i) custom_summary(blood[[i]],
19                                                         names(blood)[i]))
20 kable(summary_df, format = "latex",
21        caption = "Summary of the Predictors", booktabs = TRUE) |>
22   kable_styling(latex_options = c("striped", "scale_down")) |>
23   column_spec(1, bold = TRUE, border_right = TRUE)

```

Listing 13: R code for the Table 3

```

1 #| label: tbl-4
2 #| tbl-cap: "SD of each variable across the groups in resonse "
3 #| tbl-height: 8
4 #| tbl-width: 12
5
6 # Calculate standard deviations for all variables by disease
7 std_devs_by_disease <- blood |>
8   group_by(Disease) |>
9   summarise(across(1:23, sd, na.rm = TRUE))
10
11 # Split the table into three parts, including the first 'Disease' column in each
12 part1 <- std_devs_by_disease |> select(Disease, 1:7)
13 part2 <- std_devs_by_disease |> select(Disease, 8:14)
14 part3 <- std_devs_by_disease |> select(Disease, 15:23)
15
16 # Then, create the LaTeX tables for each part

```

```

17 kable(part1, format = "latex", booktabs = TRUE,
18     caption = "SD of Each Variable Across Disease Groups - Part 1") |>
19     kable_styling(latex_options = c("striped", "scale_down")) |>
20     column_spec(1, bold = TRUE, border_right = TRUE)
21
22 kable(part2, format = "latex", booktabs = TRUE,
23     caption = "SD of Each Variable Across Disease Groups - Part 2") |>
24     kable_styling(latex_options = c("striped", "scale_down")) |>
25     column_spec(1, bold = TRUE, border_right = TRUE)
26
27 kable(part3, format = "latex", booktabs = TRUE,
28     caption = "SD of Each Variable Across Disease Groups - Part 3") |>
29     kable_styling(latex_options = c("striped", "scale_down")) |>
30     column_spec(1, bold = TRUE, border_right = TRUE)

```

Listing 14: R code for the Table 4

References

1. Doe, J., & Smith, J. (2020). A Comprehensive Study of Blood Parameter Analysis for Disease Prediction. *Journal of Medical Research*, 10(2), 123-130.
2. Smith, J. (2021). *The Essentials of Disease Prediction in Modern Medicine*. New York: Healthcare Publishing House.
3. Aboelnaga, E. (2022). Multiple Disease Prediction Dataset. Retrieved March 22, 2024, from Kaggle