



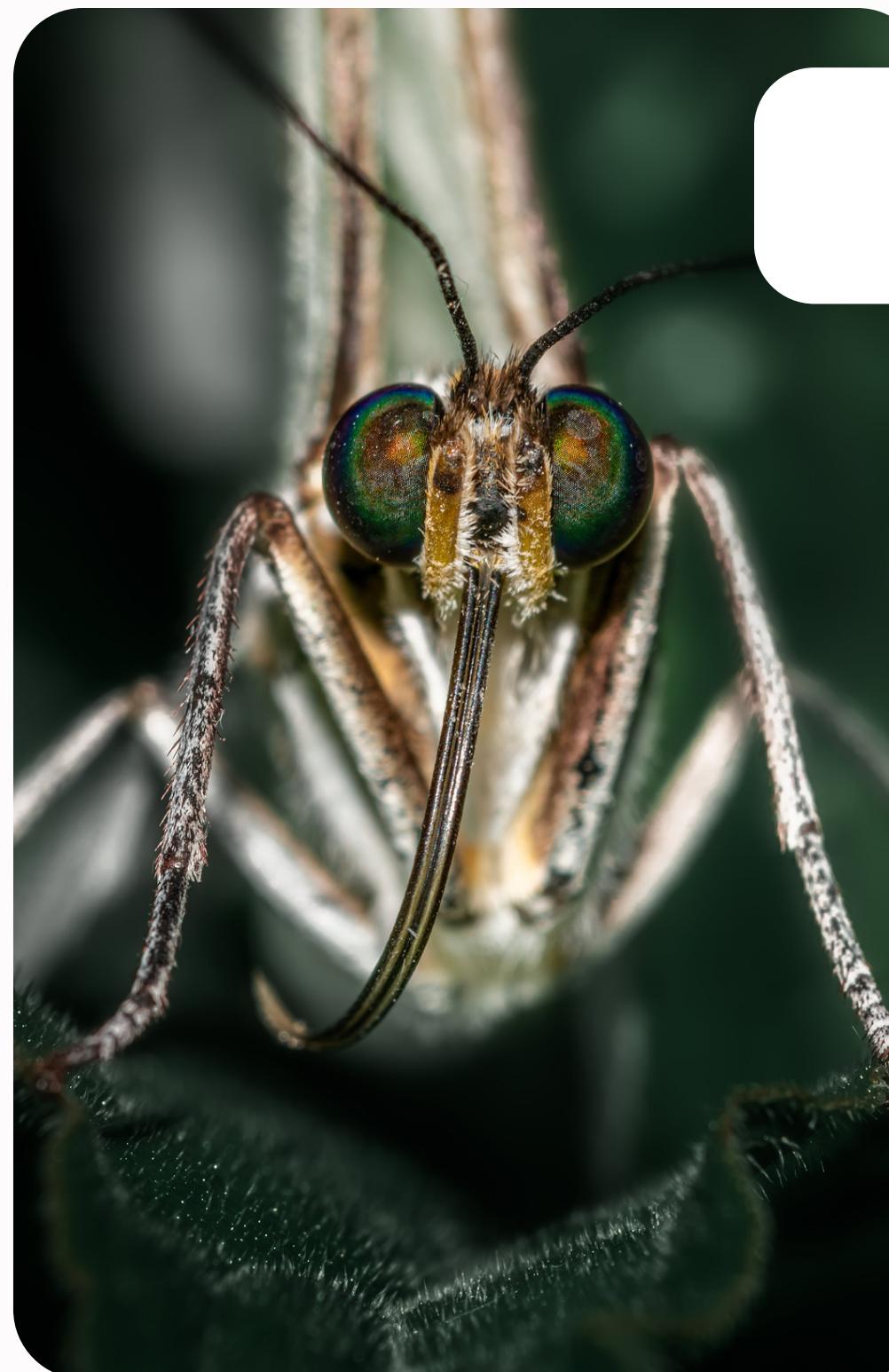
**GROUP 14**

# **CAPSTONE PROJECT**

**Malaria Prediction**



# Introduction



## Project Overview

Malaria continues to be a significant public health challenge in Kenya, with substantial efforts dedicated to its prevention and control. The Kenya Demographic and Health Survey's Malaria Indicator Survey (KDHS-MIS) provides a rich source of data that, we analyzed using advanced machine learning (ML) techniques, so as to offer valuable insights and understanding of malaria prevalence, risk factors, and intervention effectiveness.



## Problem Statement

The problem revolves around **the underutilization of machine learning techniques in analyzing Kenya's Demographic and Health Survey's Malaria Indicator Survey (KDHS-MIS) data**, hindering comprehensive understanding and predictive modeling of malaria prevalence, risk factors, and control efficacy.

Hence, Integrating machine learning into KDHS-MIS analysis aims to bridge this gap, providing actionable insights for targeted interventions and enhancing malaria control strategies in Kenya.

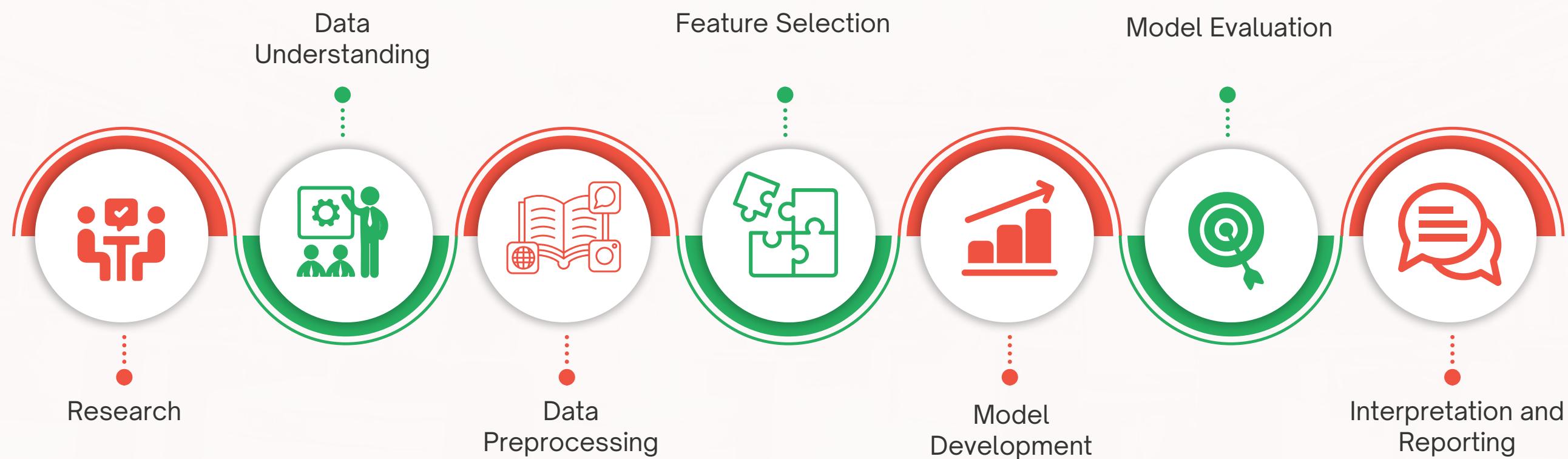
# Objectives

- O1** Apply machine learning algorithms to analyze KDHS-MIS data and identify patterns, trends, and potential predictive factors associated with malaria prevalence.
- O2** Develop predictive models for malaria occurrence, severity, and spatial distribution to aid in targeted intervention strategies.
- O3** Evaluate the effectiveness of existing malaria control measures and recommend data-driven improvements
- O4** To Conduct comprehensive feature engineering to extract relevant features from DHS data, considering demographic, socio-economic, and health-related variables.
- O5** To identify regional variations in malarial – health indicators



# Methodology

In this section, we use a methodical approach employed to analyze the Kenya Demographic and Health Survey-Malaria Indicator Survey (KDHS-MIS) dataset, that includes



# Our Approach

## Data Understanding

- Here the understanding of the data was to comprehend the nature, structure, and characteristics of the dataset.
- The data for this project was obtained from the [\[DHS Program website\]](#) ([https://dhsprogram.com/data/datasets/Kenya\\_Standard-DHS\\_2022.cfm?  
flag=0](https://dhsprogram.com/data/datasets/Kenya_Standard-DHS_2022.cfm?flag=0)).
- The data available is relevant for the intended analysis and predictions.

## Data Preprocessing

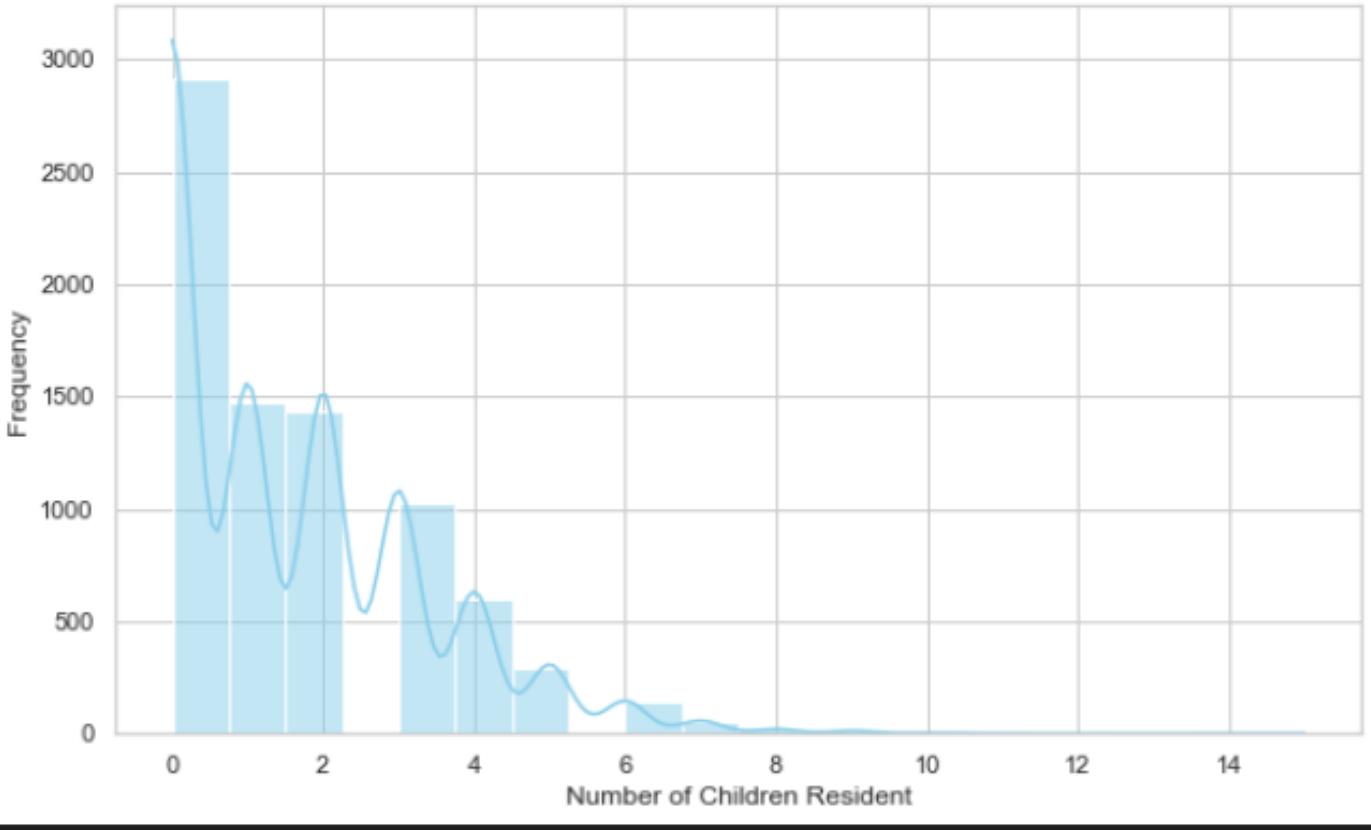
- The KDHS-MIS dataset is **cleaned and preprocessed** to handle missing values, outliers, to guarantee data quality.
- This involved:
  - 1.Collapsing repeated columns
  - 2.Handling missing values by: *Filling in missing values, Dropping empty columns*
  - 3.Handling unique values
  - 4.Identifying columns to keep
- Techniques: (simpleImputing, thresholds, fit\_transform)

## Data Visualization

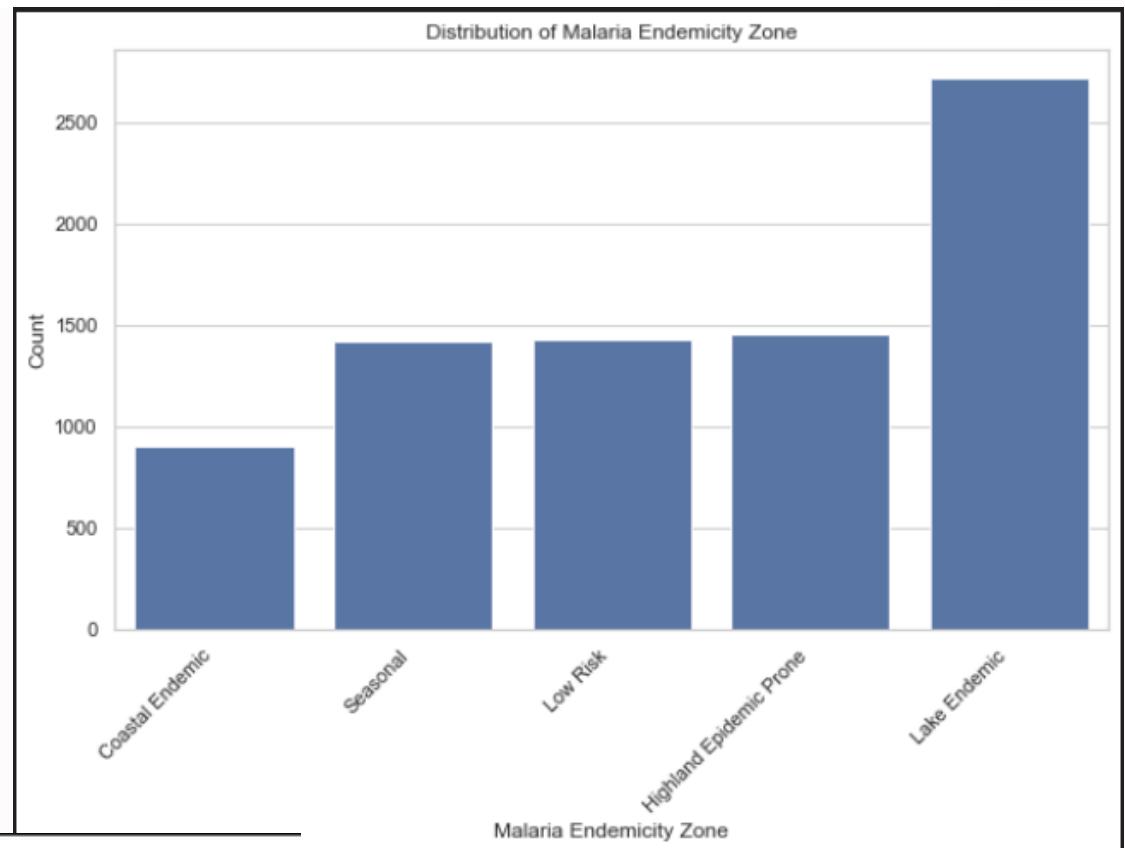
- The Exploratory Data Analysis process involved:
  1. **Univariate EDA**  
**Distribution of the various aspects that affect households** (*type\_of\_place, region\_of\_residence, water\_sources, Number of Children, Malaria Endemicity Zone*)
  2. **Multivariate EDA**  
**Relationship of aspects that affect households** (*correlation matrix, type\_of\_place Vs Availability of Mosquito nets, malaria cases by region, malaria cases by age group*)

# Univariate EDA

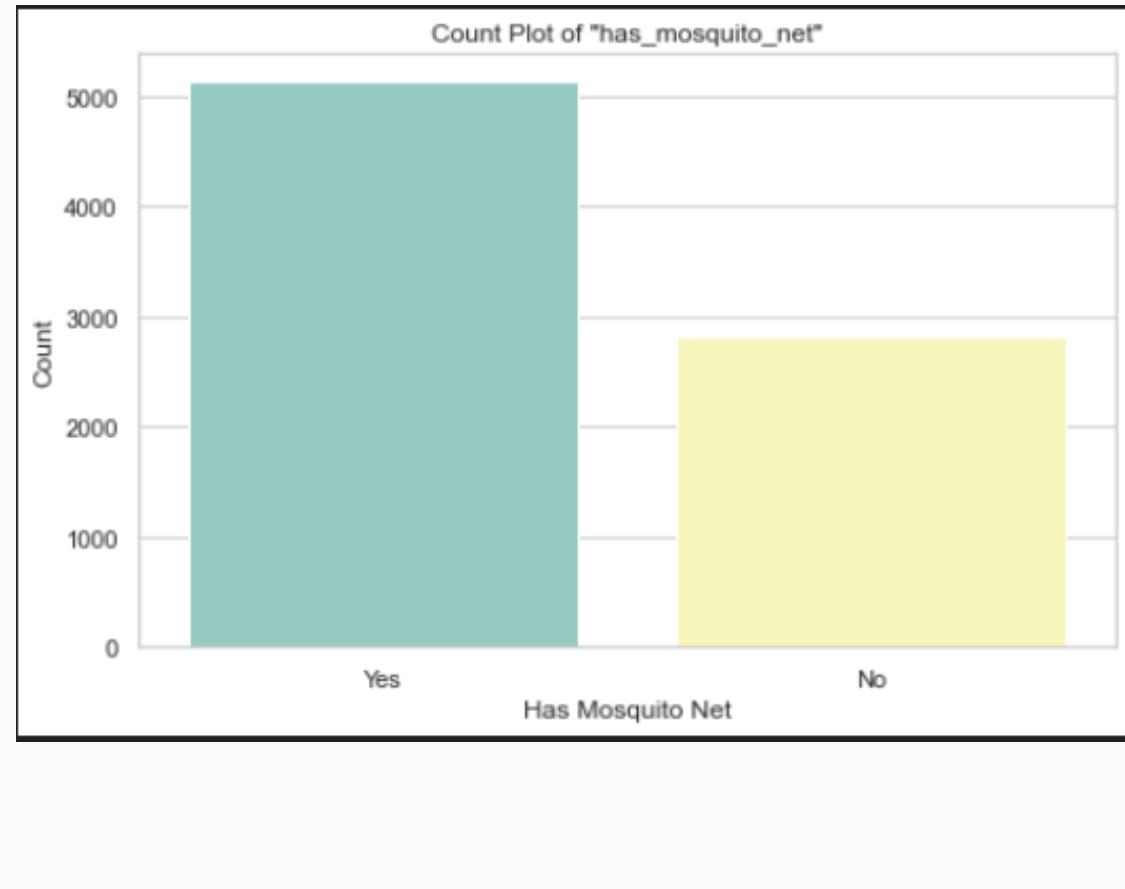
Distribution of Number of Children Resident in Each Household



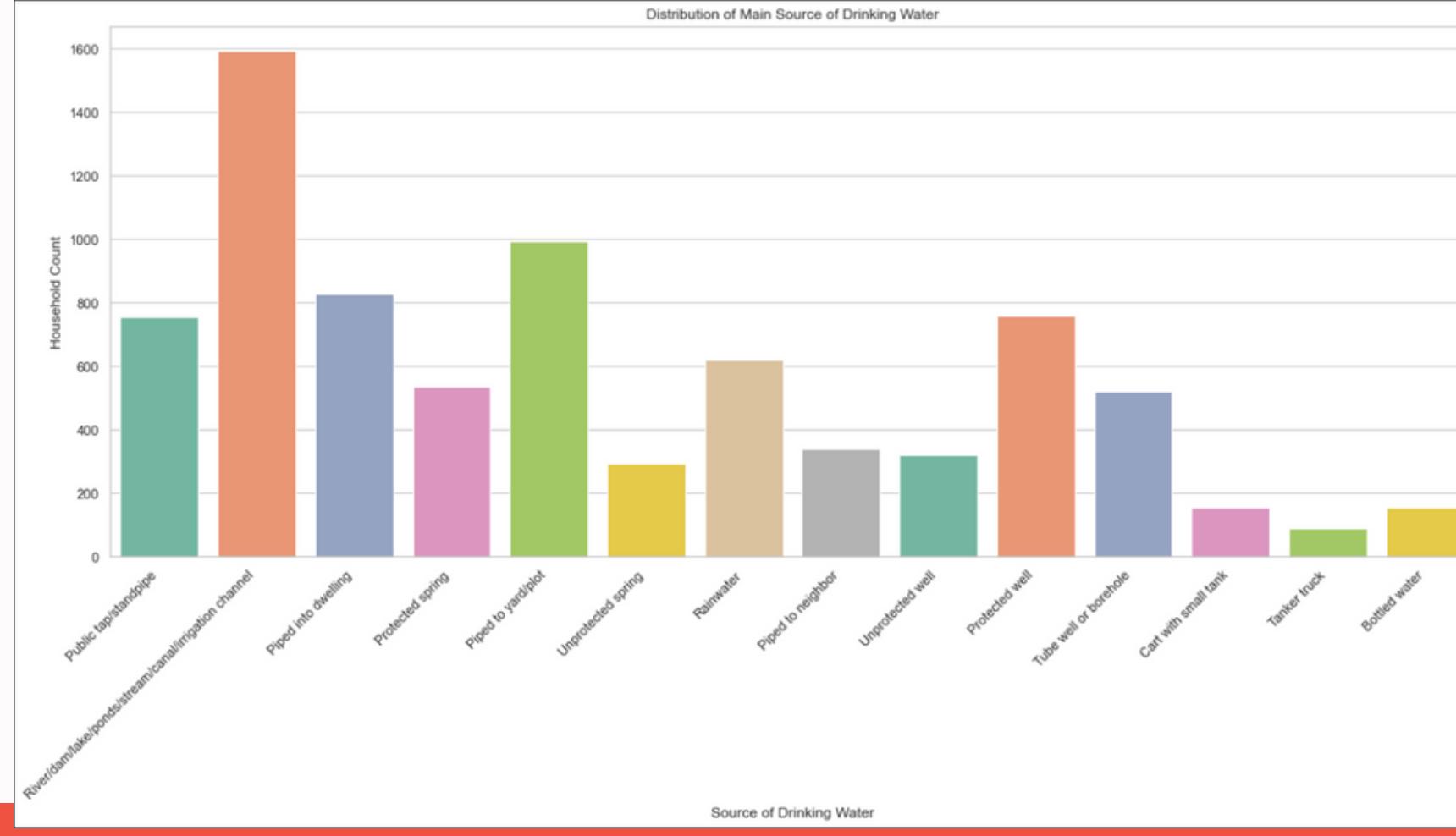
Distribution of Malaria Endemicity Zone



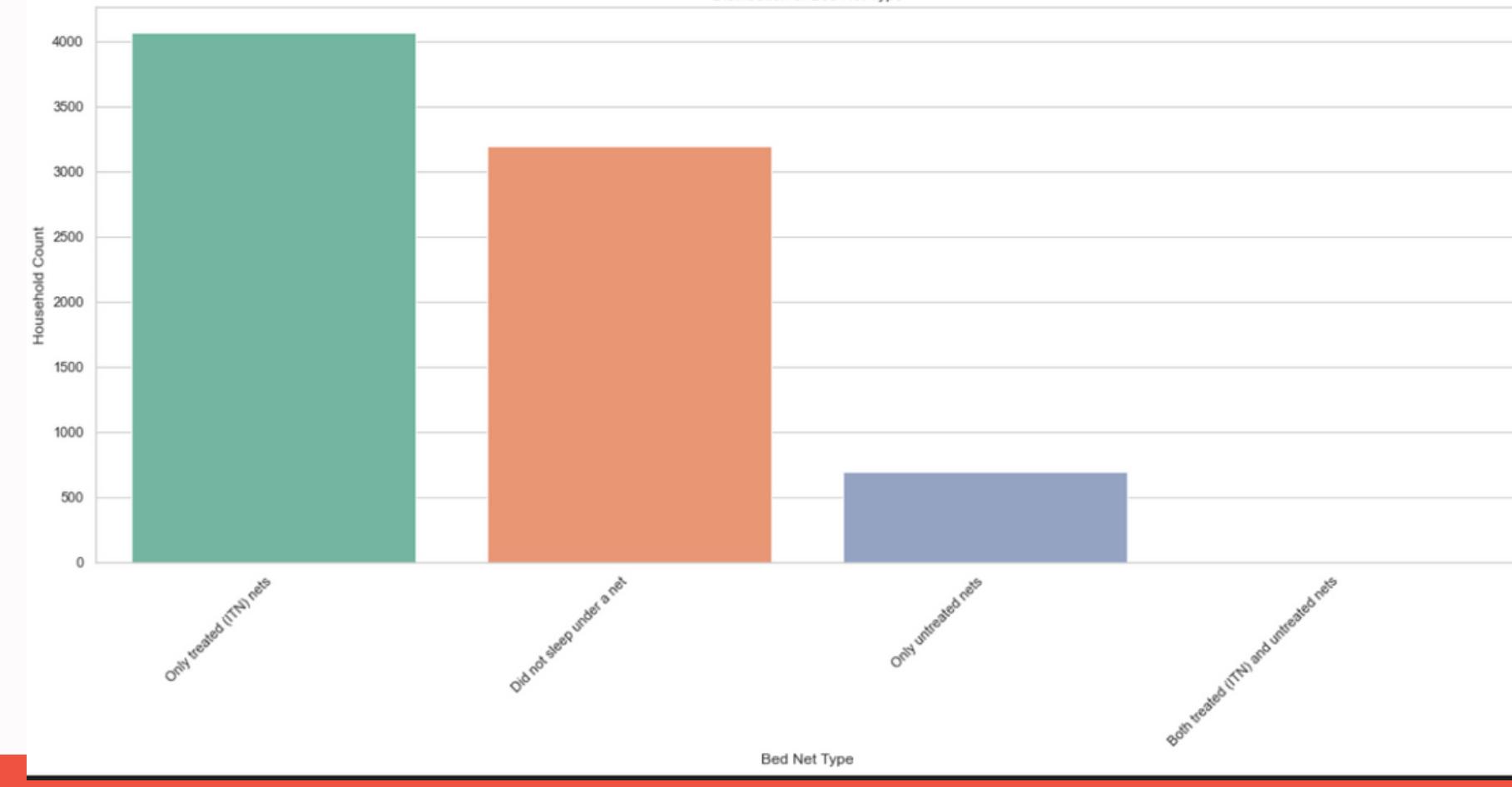
Count Plot of "has\_mosquito\_net"



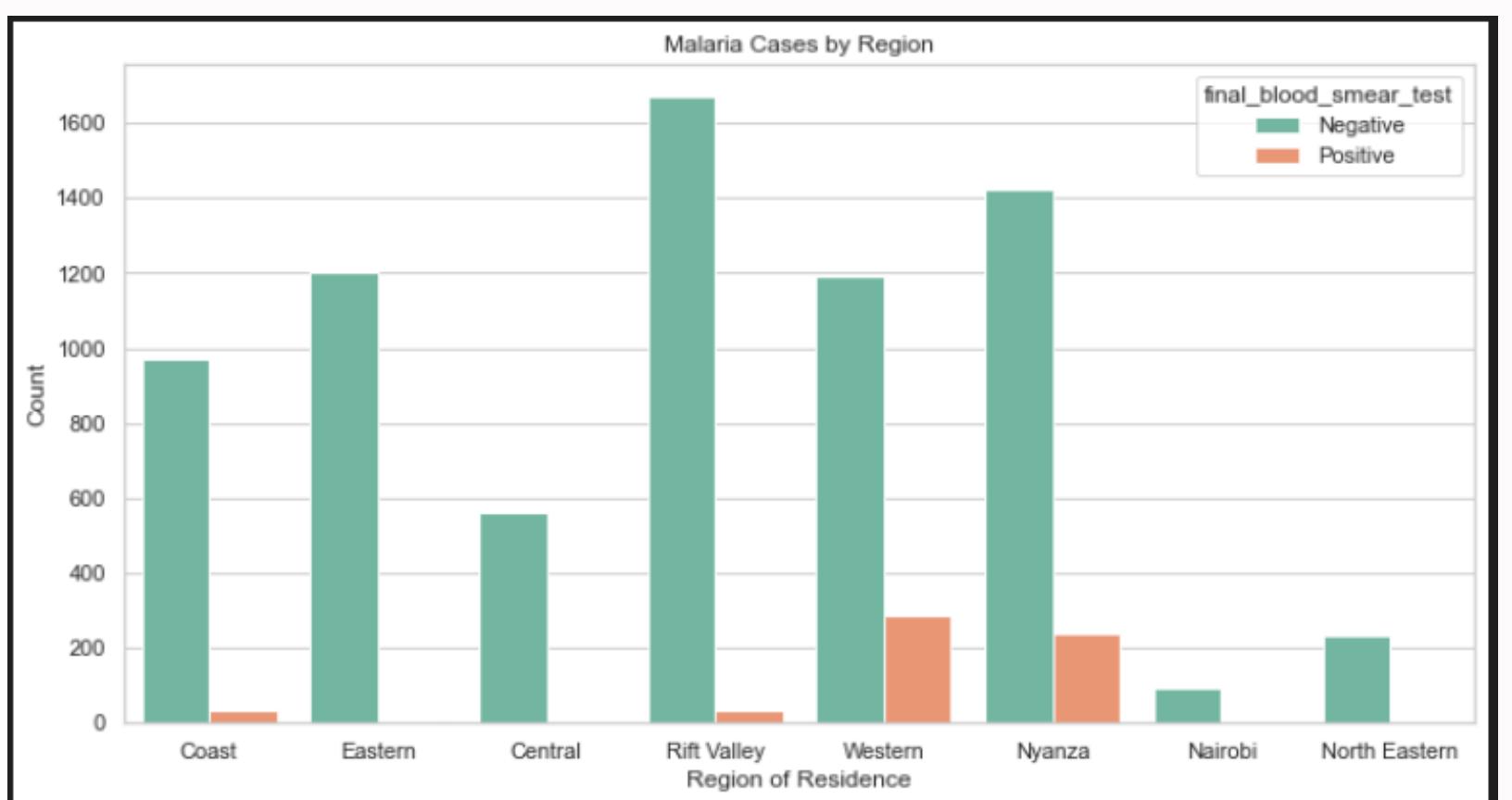
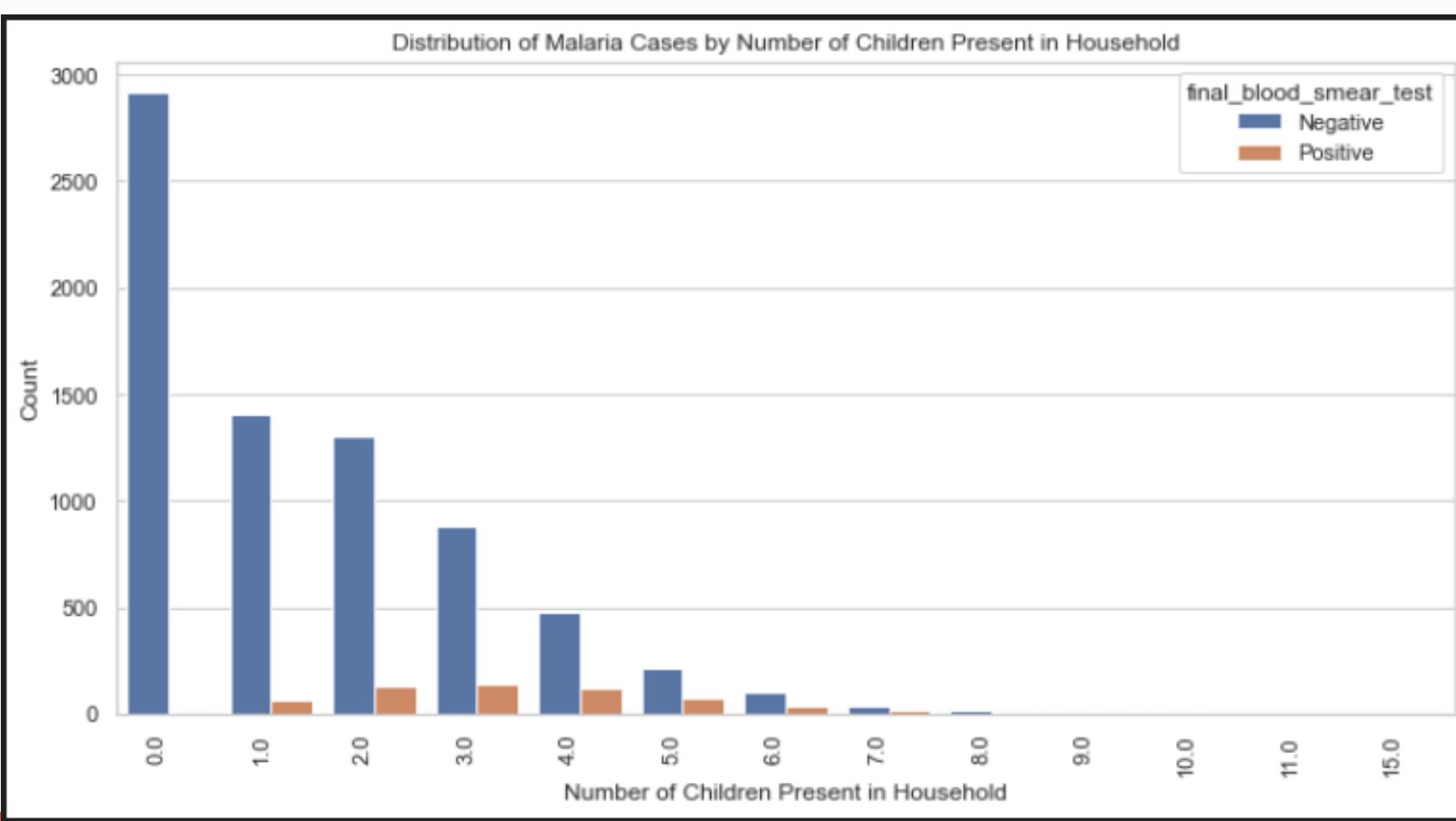
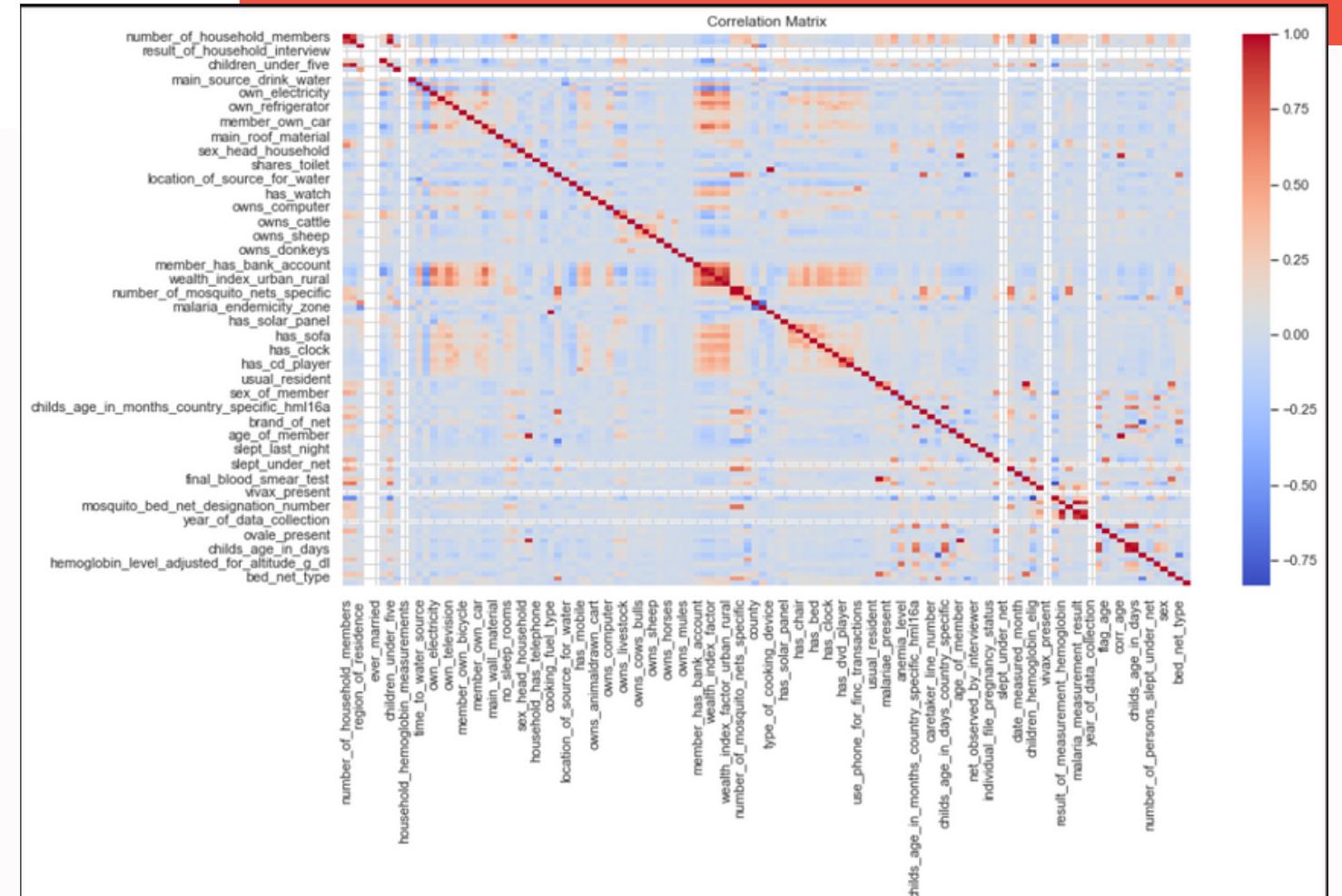
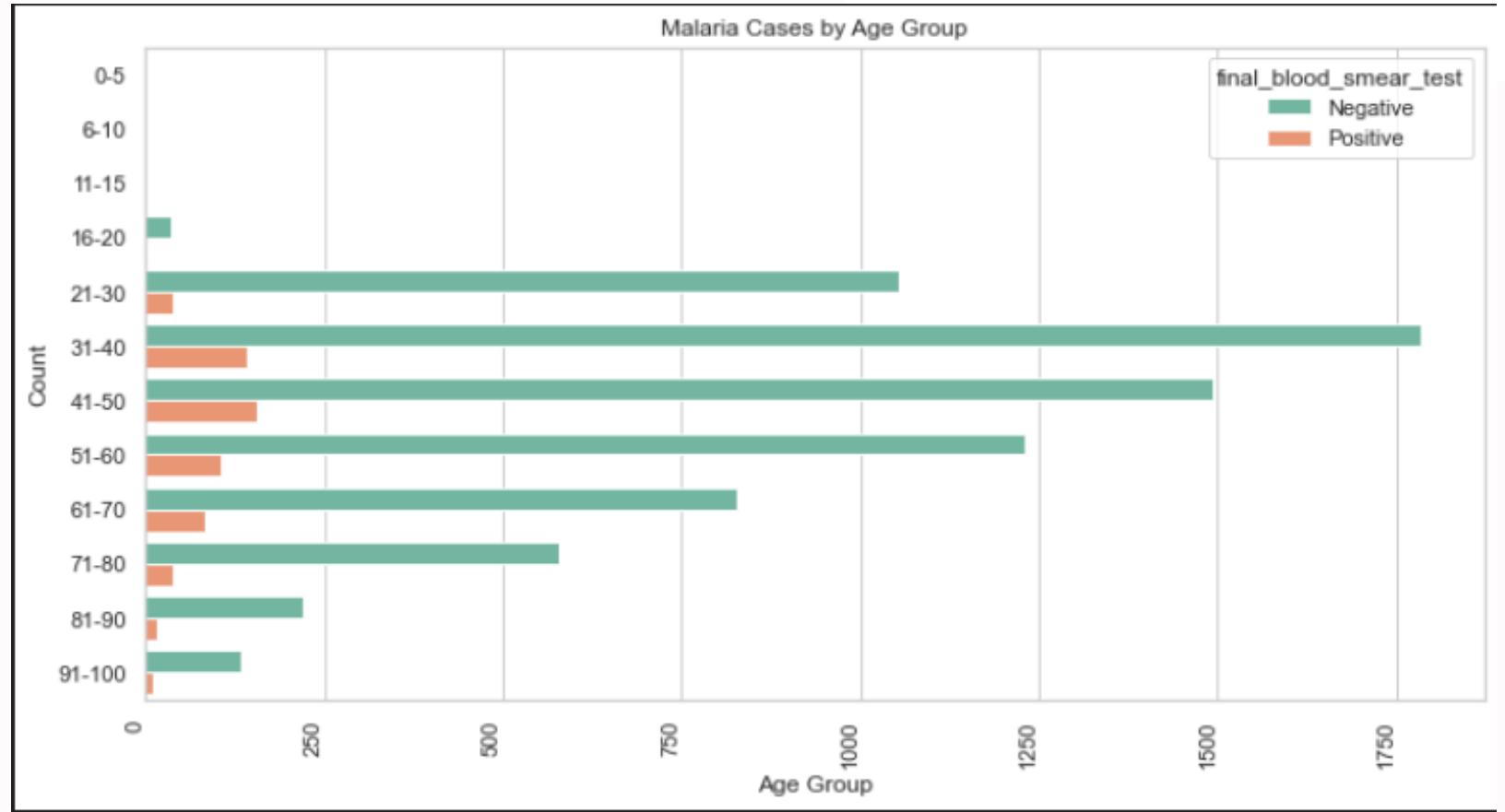
Distribution of Main Source of Drinking Water



Distribution of Bed Net Type



# Multivariate EDA



# Modelling

## Logistic Regression model

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam non tincidunt enim. Sed rutrum lorem eu bibendum pretium. Duis bibendum semper.

## Lasso Model

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam non tincidunt enim. Sed rutrum lorem eu bibendum pretium. Duis bibendum semper.

## Random Forest model

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam non tincidunt enim. Sed rutrum lorem eu bibendum pretium. Duis bibendum semper.



# Model Evaluation

An interpretation and explanation of the key metrics:

## Accuracy

Accuracy measures the overall correctness of the model predictions.

In this case, the model has an accuracy of approximately 96.67%, indicating that about 96.67% of the predictions are correct.

The model is performing well overall with high accuracy.

The precision for positive results is relatively high, indicating that when the model predicts a positive result, it is likely to be correct.

The recall for positive results is also reasonable, suggesting that the model is effective at identifying a significant portion of the actual positive instances.

The F1-score reflects a good balance between precision and recall.

The model seems to be particularly effective at correctly predicting negative results (0.0), as indicated by the high precision and recall for this class.

## Precision

Precision measures the accuracy of the positive predictions made by the model.

Precision for positive results (1.0) is approximately 90%. It means that when the model predicts a positive result, it is correct about 90% of the time.

## Recall (sensitivity)

Recall measures the model's ability to correctly identify positive instances.

Recall for positive results (1.0) is approximately 72%. It means that the model correctly identifies about 72% of the actual positive instances.



# REPORTING

## Recommendation

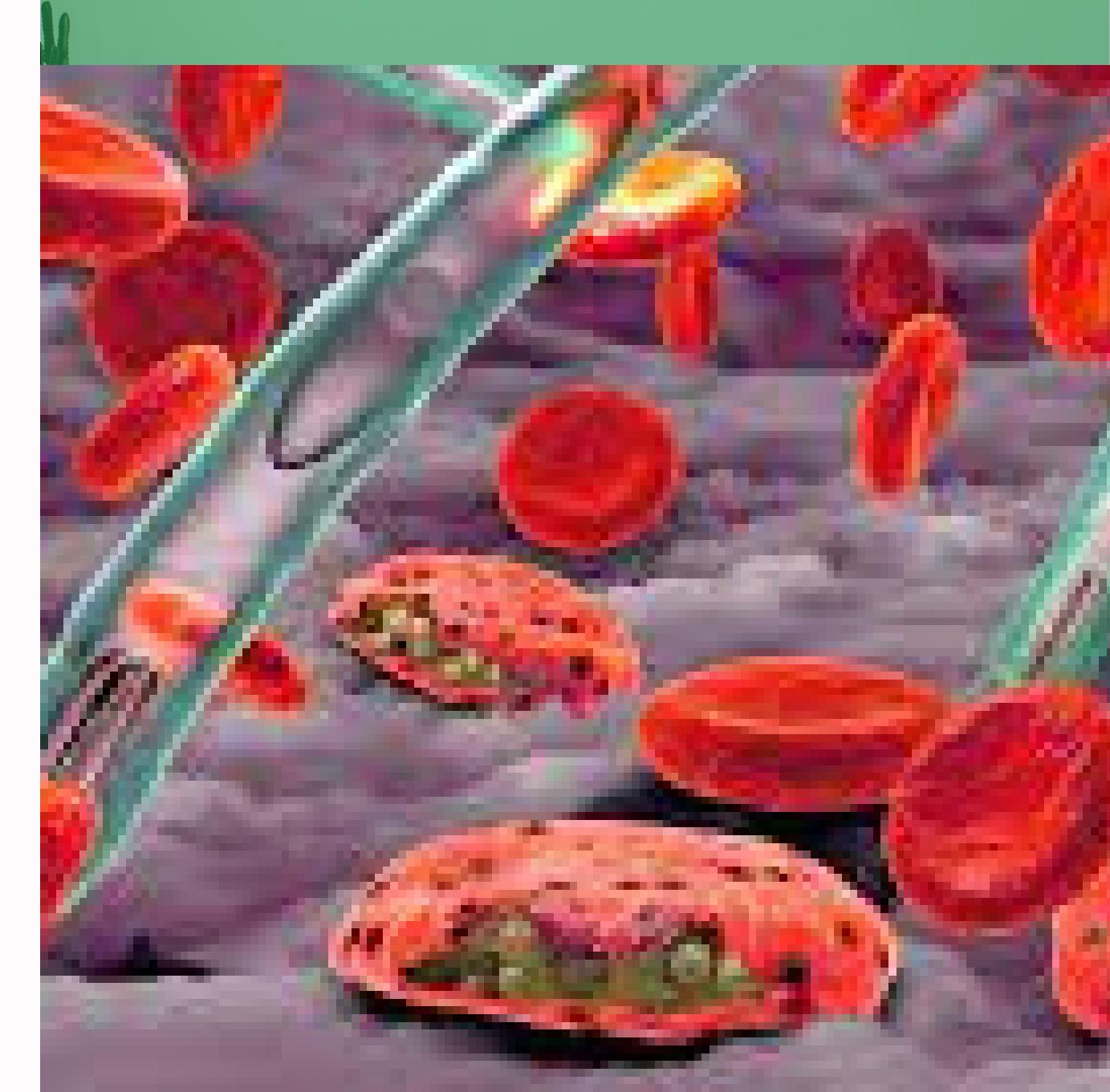
The model does accurate predicting of malaria rapid test results.

MCC value of 0.7833 indicates strong overall model performance, taking into account both false positives and false negatives. This metric is particularly useful in situations where there is class imbalance or where the consequences of different types of errors vary.

## Conclusion

The model does accurate predicting of malaria rapid test results.

MCC value of 0.7833 indicates strong overall model performance, taking into account both false positives and false negatives. This metric is particularly useful in situations where there is class imbalance or where the consequences of different types of errors vary.



# DEPLOYMENT

- Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam non tincidunt enim. Sed rutrum lorem eu bibendum pretium. Duis bibendum semper. Tincidunt enim bibendum pretium sit amet. Consectetur adipiscing elit.
- Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam non tincidunt enim. Sed rutrum lorem eu bibendum pretium. Duis bibendum semper. Tincidunt enim bibendum pretium sit amet.



# Future Work

- To visualize regional variations in malaria risk indicators using; heatmaps, gmap
- To find the relationship between the malaria indicators and most common diseases in the different demographics.
- Analyze trends in malaria prevalence over time.



## OUR TEAM

- 
1. Alpha Guya
  2. Ben Ochoro
  3. Caleb Ochieng
  4. Christine Mukiri
  5. Dominic Muli
  6. Frank Mandele
  7. Jacquline Tulinye
  8. Lesley Wanjiku



**THANK  
YOU**