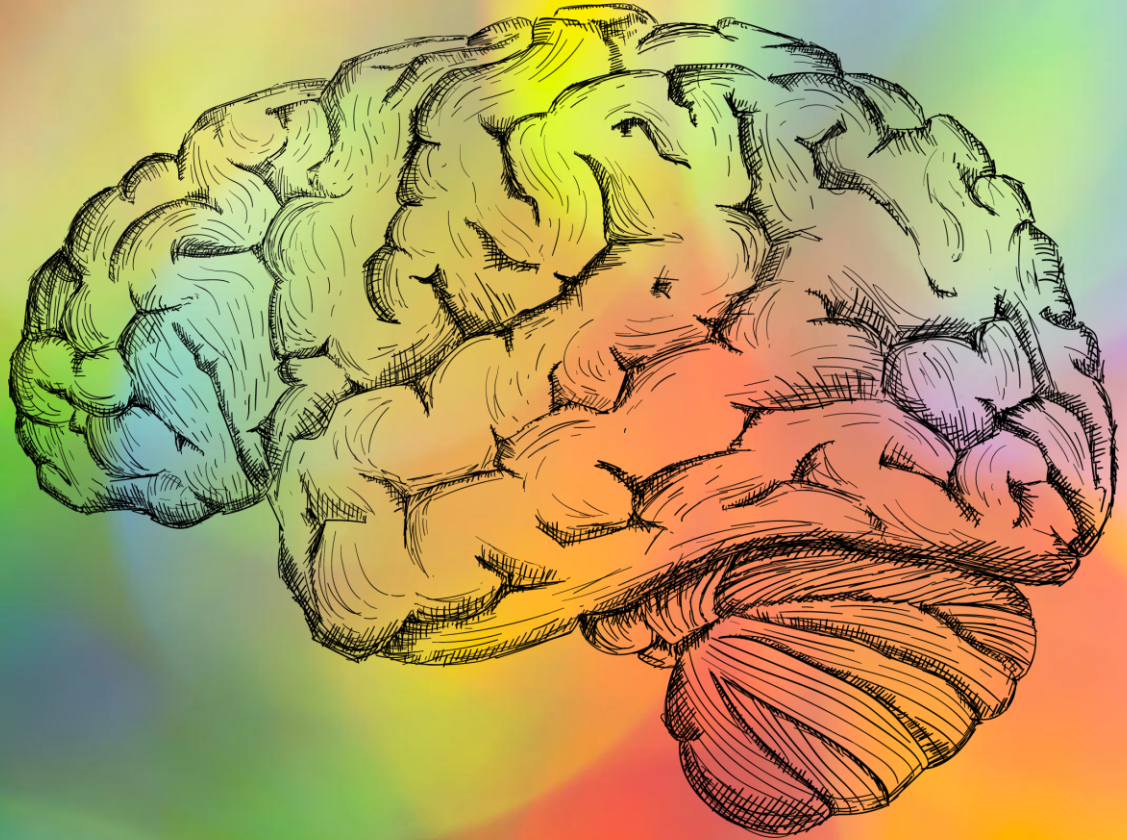# Using Data Science to improve Stroke Prediction

By:

Josiah Joel – 816030501

Ethan Lee Chong – 816032732

# Research Goal

**To compare the Kappa algorithm with other traditional data science algorithms using a real-world stroke prediction data set and a synthetic dataset.**

# - Kappa deals with numeric columns differently

Train set:

- if the mark 26.2 appears in the testing
the algorithm wasn't trained with that value
so....
what does the algorithm do?????

| categoric column | NUMERIC COLUMN | categoric column |
|---|---|---|
| Sex | Mark | Pass |
| F | 29.3 | No |
| M | 31.4 | No |
| O | 50.0 | Yes |
| F | 51.7 | Yes |

puts all numbers
on one scale
so it's easy to guess

# What is a stroke?

• **Strokes are the second most common cause of <u>death</u>, and it is one of the leading causes of disability worldwide. (Cleveland Clinic, 2022)**

• **A stroke can be caused by a blockage of blood supply to the brain or when a blood vessel in the brain bursts. (CDC, 2023) There are two main types of stroke: Ischemic and Haemorrhagic.**

# Types of stroke

## Ischemic

- caused by plaque or fatty deposit in blood vessels near or in the brain

- Ischemic strokes are the more common of the two
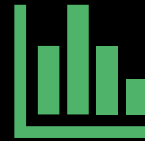
## Haemorrhagic

- caused by brain aneurysms, brain tumors, blood thinning medications, head injuries and Ischemic stroke that had secondary bleeding.

Stroke chance is mainly affected by lifestyle choices.

# Beneficiaries

Helps medical industry with accurate predictions.

General public benefits from analytics and predictions.

Predict high-risk individuals for timely intervention.

The Kappa algorithm compares each row of test data with all the training data. This maximizes the robustness.

The weights are also calculated using the numerical and categorical means maximizing the personalization.

$$\hat{y}_j[k] \equiv \frac{\sum_{i \in S} \frac{y_i}{(1+|x_i-x_j|)^{\kappa^*}}}{\sum_{i \in S} \frac{1}{(1+|x_i-x_j|)^{\kappa^*}}}$$

Don't worry. This is the important part.
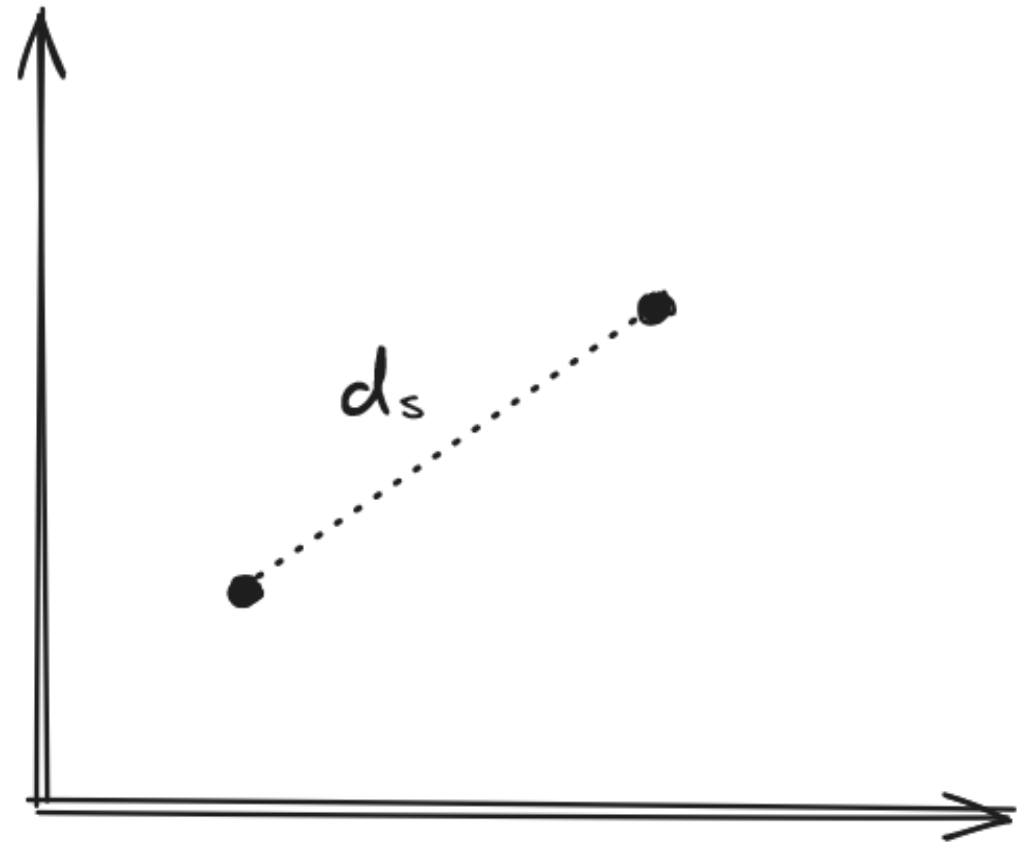Kappa constant.

Same formula but less scary:

$$g(\kappa) = \frac{\sum_{s \in S} w(s, \kappa) M(s)}{\sum_{s \in S} w(s, \kappa)}$$

Let's simplify this formula:

Weighted average w:

$$w(s, \kappa) = \frac{1}{(1 + d_s)^{\kappa}}$$

$d_s$

# Significance of the Project Solution

If the algorithm proves to be better at classification in this dataset than traditional approaches it could be applied to other real-world problems and problems with many more features. Many lives can also be saved from improvements.

# Goals and Objectives

- Research paper.

- Github repo with notebook files.

- F1 scores, Accuracy scores and Weighted Incorrect Classification

Metric of all algorithms used.

# F1 scores

# F1 scores

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

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives}+\text{False Negatives}}$$

$$F1 = \frac{2\times\text{Precision}\times\text{Recall}}{\text{Precision}+\text{Recall}}$$

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | TN | FP |
| Actual 1 | FN | TP |

# Accuracy scores

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{Total Population}}$$

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | TN | FP |
| Actual 1 | FN | TP |

# WICM scores

Weighted Incorrect Classification Metric =

$$\frac{Weight \times False\ Negatives + False\ Positives}{True\ Positives + True\ Negatives + False\ Positives + False\ Negatives}$$

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | TN | FP |
| Actual 1 | FN | TP |

We set our weight arbitrarily to 10

F1 Scores for Different Machine Learning Models

# Our Dataset is very skewed

Around **95%** of the rows in the dataset are negative for stroke. We are predicting stroke, and this is not good. To solve this issue, we use:
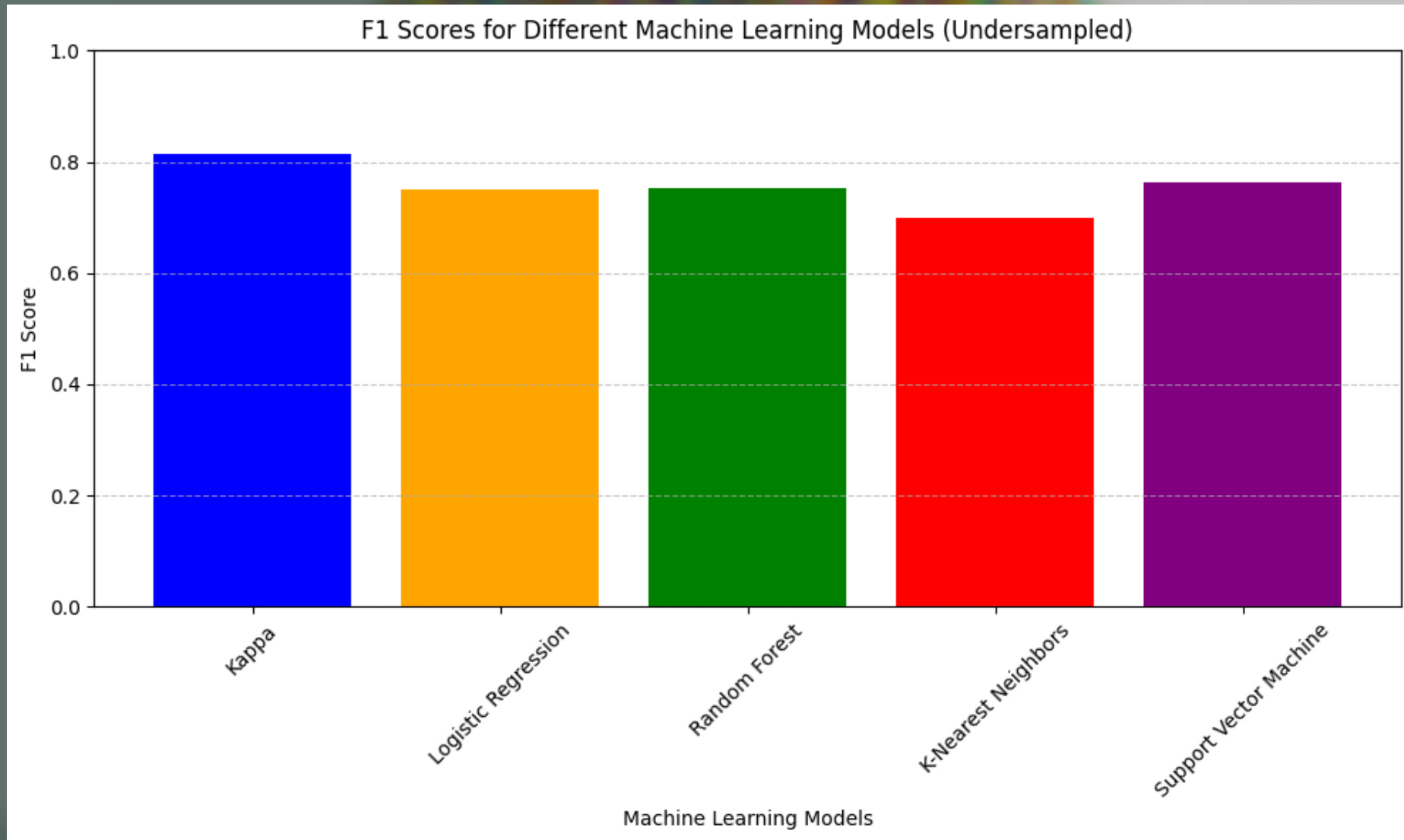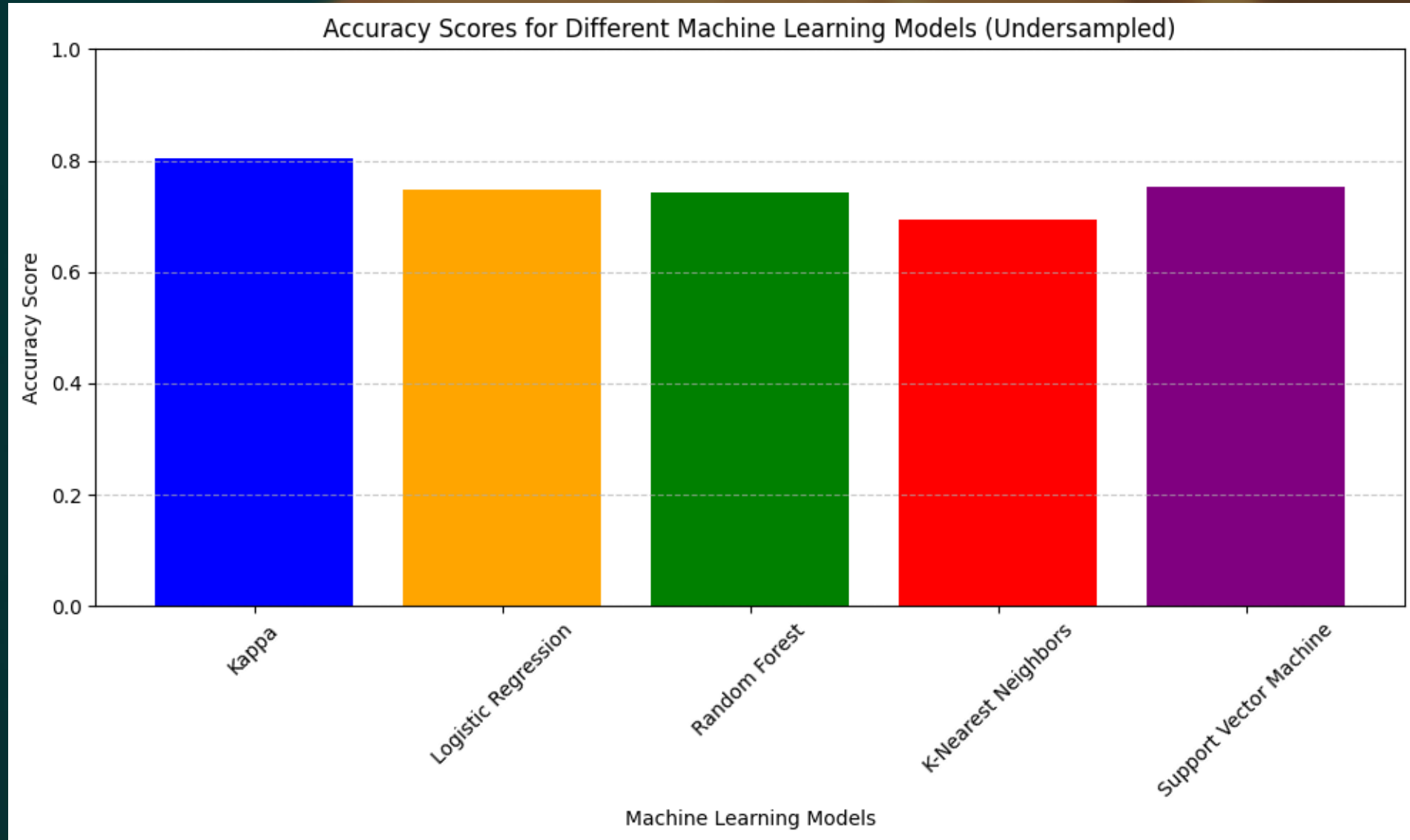
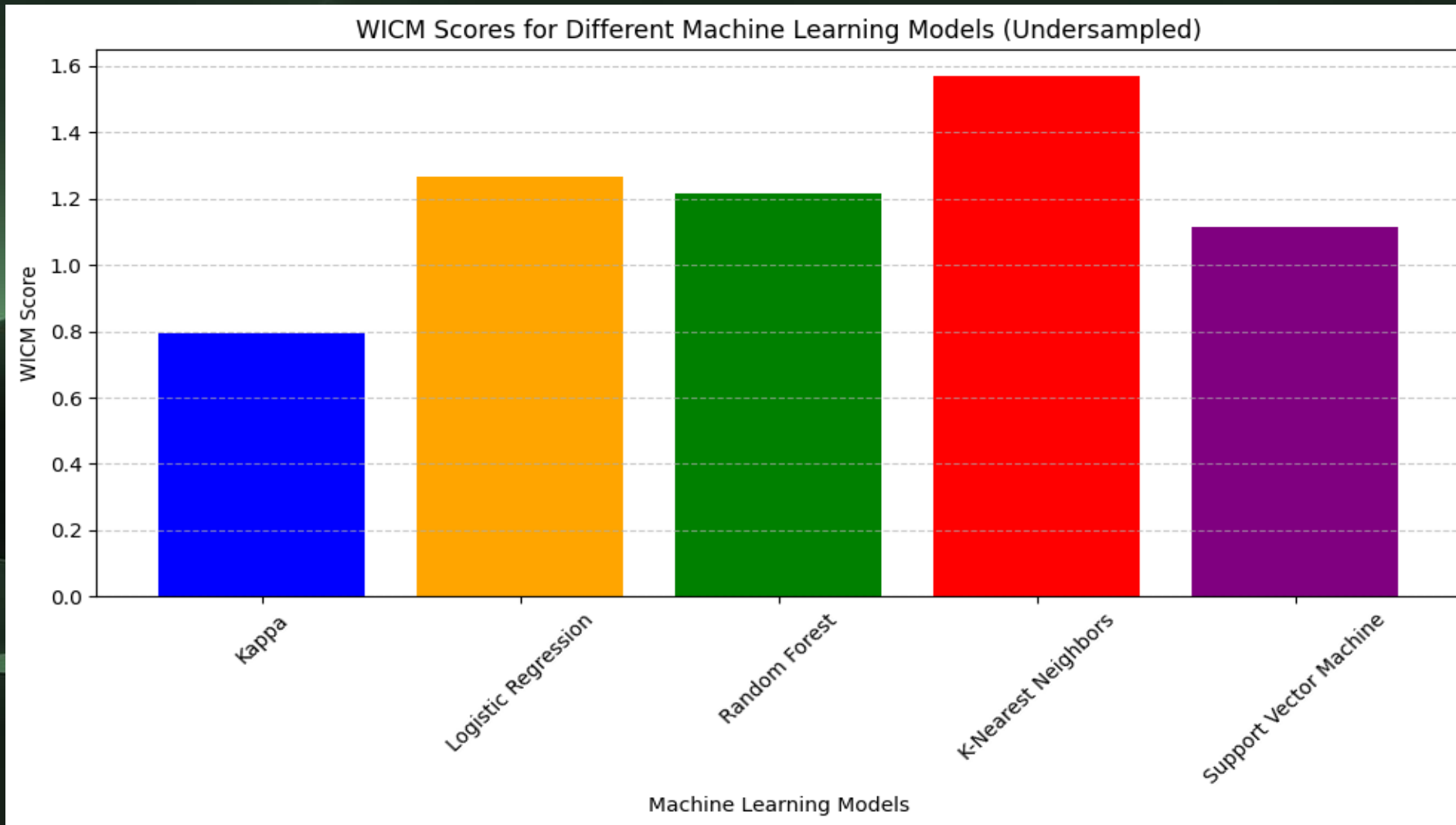## SMOTE OVERSAMPLING

and

## UNDERSAMPLING

# Undersampled F1 scores



F1 Scores for Different Machine Learning Models (Undersampled)

# Undersampled Accuracy scores



Accuracy Scores for Different Machine Learning Models (Undersampled)

# Undersampled WICM scores



WICM Scores for Different Machine Learning Models (Undersampled)

Lower is better

# Oversampled F1 scores



F1 Scores for Different Machine Learning Models (Oversampled)

# Oversampled F1 scores



Accuracy Scores for Different Machine Learning Models (Oversampled)

# Oversampled F1 scores



WICM Scores for Different Machine Learning Models (Oversampled)
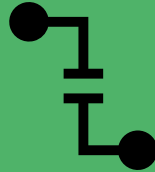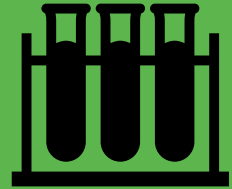
# Future Works

**TEST WITH MORE VARIED DATASETS**

**OPTIMIZE TO FIND IDEAL HYPER-PARAMETERS**

**CONDUCT FURTHER EXPERIMENTS**

# Conclusion

- Kappa balances personalization and robustness in insurance premiums

- Research on stroke prediction- Aims for more accurate risk assessments

- Potential to revolutionize predictive analytics

- Evaluated via F1 score, Accuracy and WICM

- Transparent documentation on GitHub

- Goals: advance healthcare and risk management applications and reduce the risk of death