

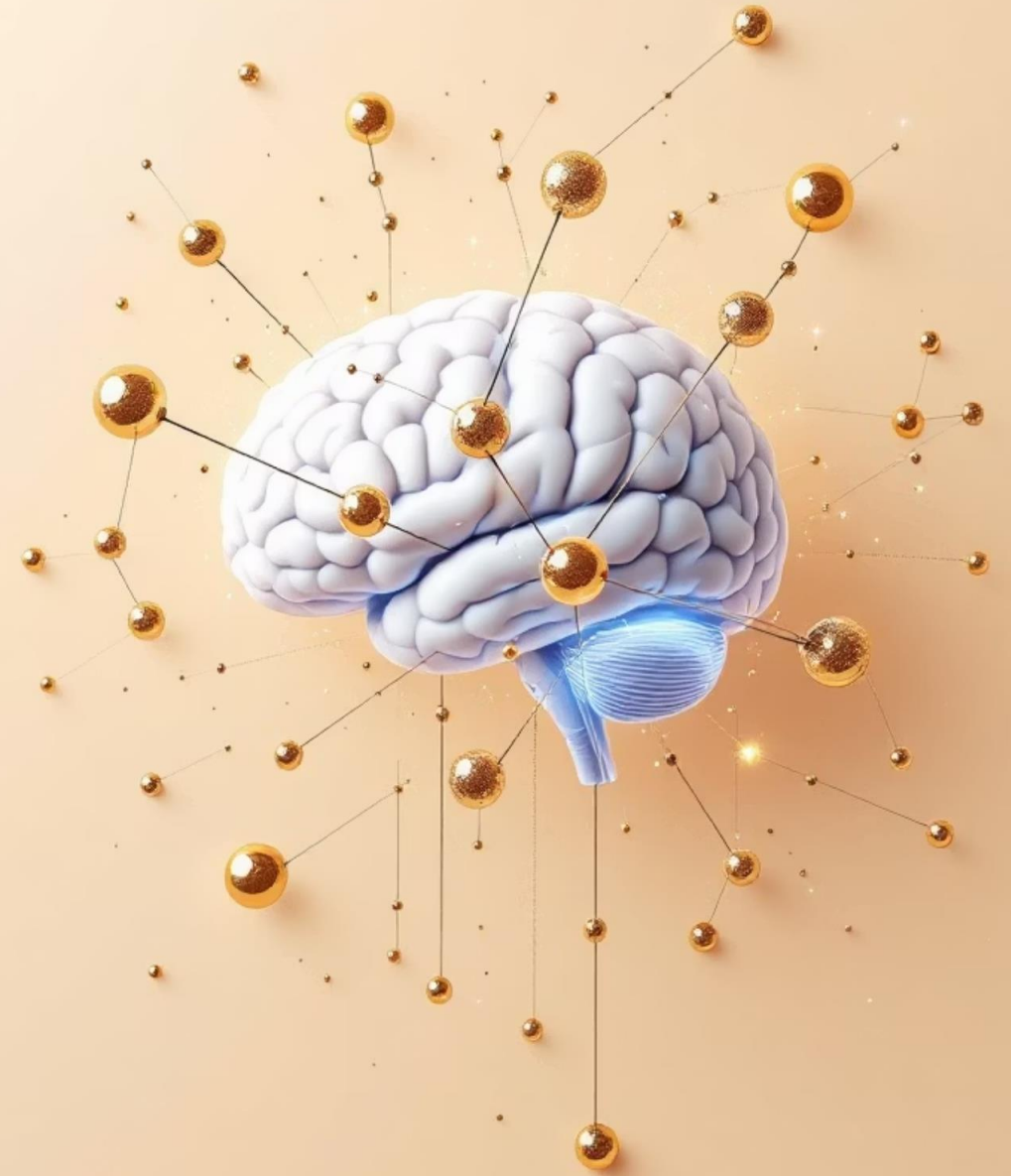


ML Classification for Stroke Risk Prediction

A Comprehensive Machine Learning Process for Healthcare Analytics

Presented by Natnael Abayneh

December 2024



Problem Definition – The Healthcare Context

The Business Problem

Stroke is a critical health emergency where every minute matters.

Research Motivation

Identifying high-risk individuals through clinical data can significantly reduce mortality and long-term disability rates.

Goal

Develop a binary classification model to predict "Stroke" (1) or "No Stroke" (0) based on patient health indicators.

Project Goals & Objectives

Primary Objective

Build a high-accuracy classification system.

Focus Areas

Focus on early detection using non-invasive features like Age, BMI, and hypertension status.

Scope

This project covers the full end-to-end ML lifecycle from raw data to a production-ready model.

Defining Success Criteria

Accuracy Target

Achieve a validation accuracy of over 85%.

Reliability Metrics

Focus on High Precision (minimizing false alarms) and High Recall (ensuring no stroke is missed).

Generalization Goal

Maintain a training-to-validation accuracy gap of less than 0.1 to avoid overfitting.

Data Collection & Sources



Source

Publicly available Healthcare
Stroke Prediction Dataset.



Dataset Size

5,110 patient records.



Domain Representation

The data includes diverse
features covering
demographics, clinical history,
and lifestyle factors.

Feature Breakdown – Feature Dictionary

Categorical Features

- Gender
- Ever Married
- Work Type
- Residence Type
- Smoking Status

Numerical Features

- Age
- Hypertension (0/1)
- Heart Disease (0/1)
- Average Glucose Level
- BMI

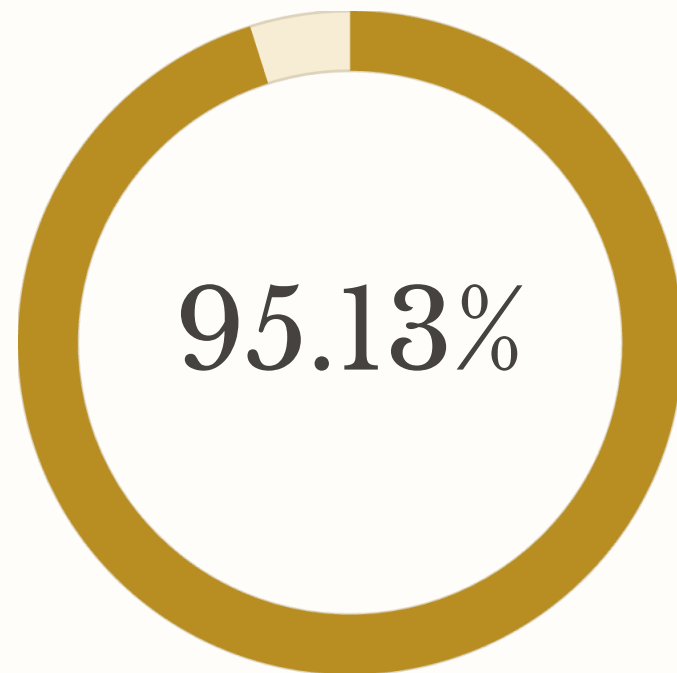
Target

stroke (0 = No, 1 = Yes)



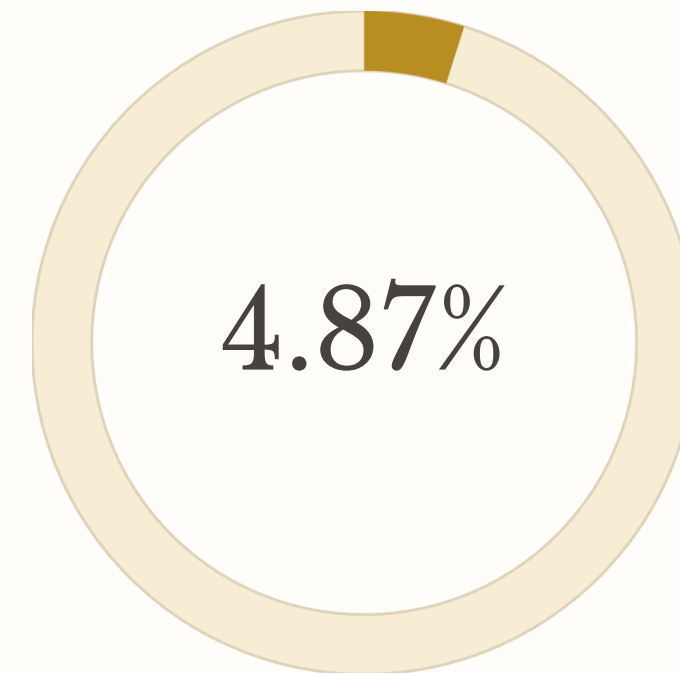
EDA – Target Variable Distribution

The dataset is "Imbalanced".



No Stroke

Majority Class



Stroke

Minority Class

Roughly 4.87% of patients had a stroke, while 95.13% did not.

The model must be carefully trained to recognize the "minority" class (actual strokes).

EDA – Identifying Key Feature Relationships

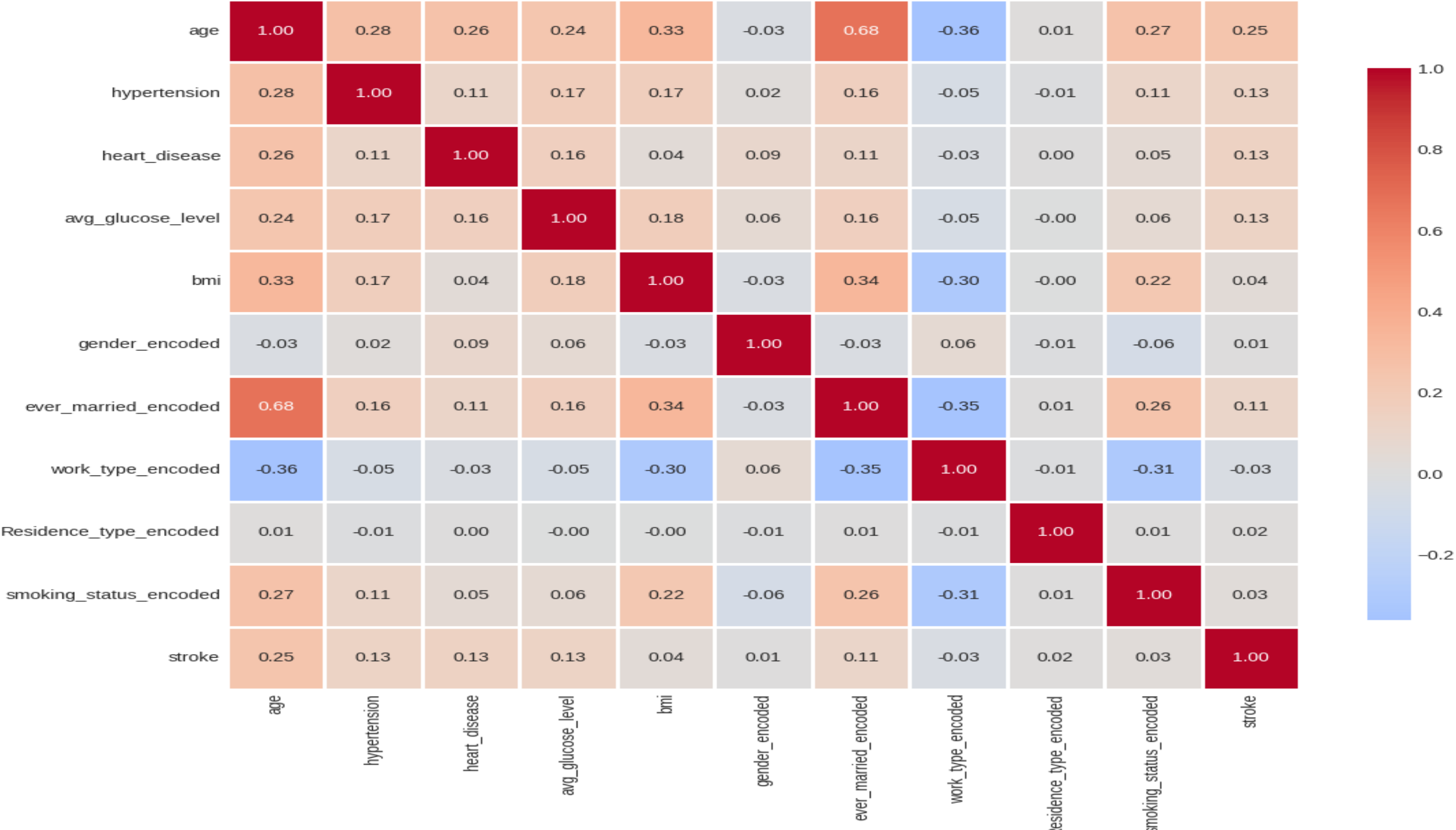
Correlation Analysis

Exploring how age and BMI correlate with stroke occurrence.

Trends

Higher age and the presence of hypertension show the strongest positive correlation with stroke risk.

Correlation Matrix



Data Cleaning – Handling Missing Values



1

Detection

2

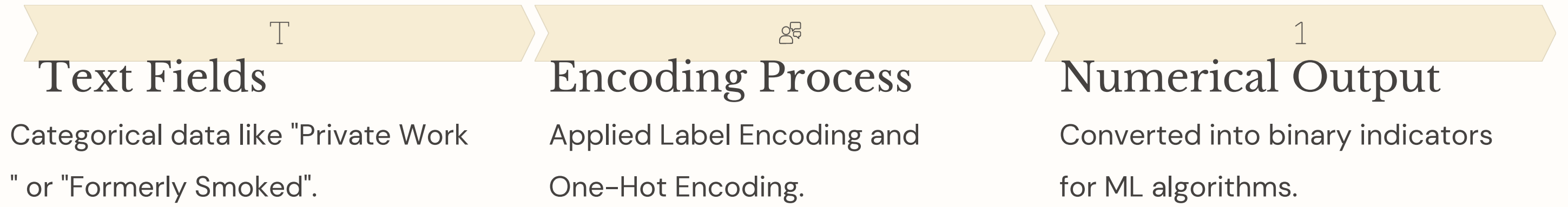
Solution

3

Integrity

Feature Engineering – Categorical Encoding

Machine learning algorithms require numerical input.



This process converts text fields (e.g., "Private Work" or "Formerly Smoked") into binary indicators, enabling machine learning algorithms to process the data effectively.

Feature Engineering – Feature Scaling

The Challenge: Disparate Scales

Features like Age (up to 82) and Glucose Level (up to 271) have different numerical scales, meaning some features inherently have larger values than others.

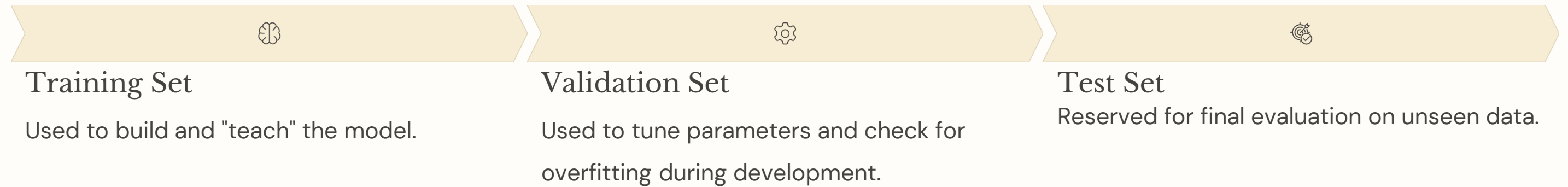
Without addressing this, models might implicitly give more weight to features with larger numerical ranges, causing them to dominate the learning process and potentially leading to suboptimal or biased model performance.

The Solution: Standard Scaling / Normalization

Standard Scaling and Normalization techniques transform features so they are all on a similar scale. Normalization typically scales values to a range between 0 and 1, while Standard Scaling transforms data to have a mean of 0 and a standard deviation of 1.

This ensures that no single feature dominates the model purely based on its magnitude, allowing the algorithm to learn genuine relationships and contributions from all features effectively.

Data Splitting Strategy



We split the data into Training, Validation, and Test sets to ensure robust model development and evaluation.

This systematic approach helps in developing a model that generalizes well to new, unseen data and avoids common pitfalls like overfitting.

Algorithm Selection Rationale

Models Evaluated

Logistic Regression, Random Forest, Support Vector Machine (SVM), and Neural Networks.

Selection Criteria

Chose algorithms based on their ability to handle imbalanced data and their computational efficiency.

Model Development & Training Workflow



Architecture

Designed a robust classification pipeline.



Training Process

The model was trained iteratively, adjusting weights to improve its detection of stroke cases.

Hyperparameter Tuning Strategy



Optimization

Used GridSearchCV and RandomizedSearchCV to find the best internal settings for the algorithms.



Fine-tuning

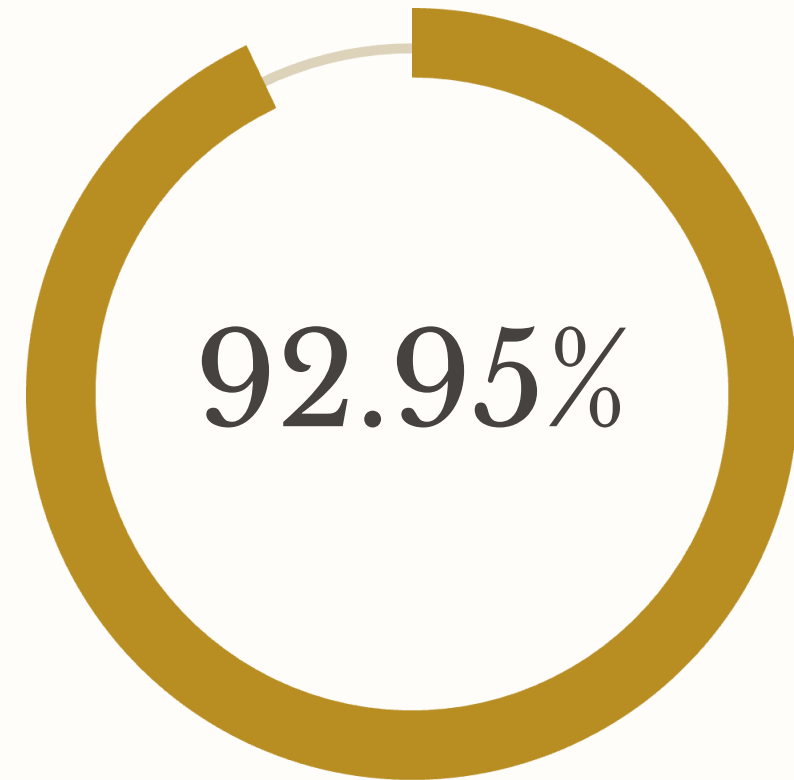
Adjusted parameters like tree depth and learning rates to maximize performance.

Classification Results & Metrics



Training Accuracy

99.58%

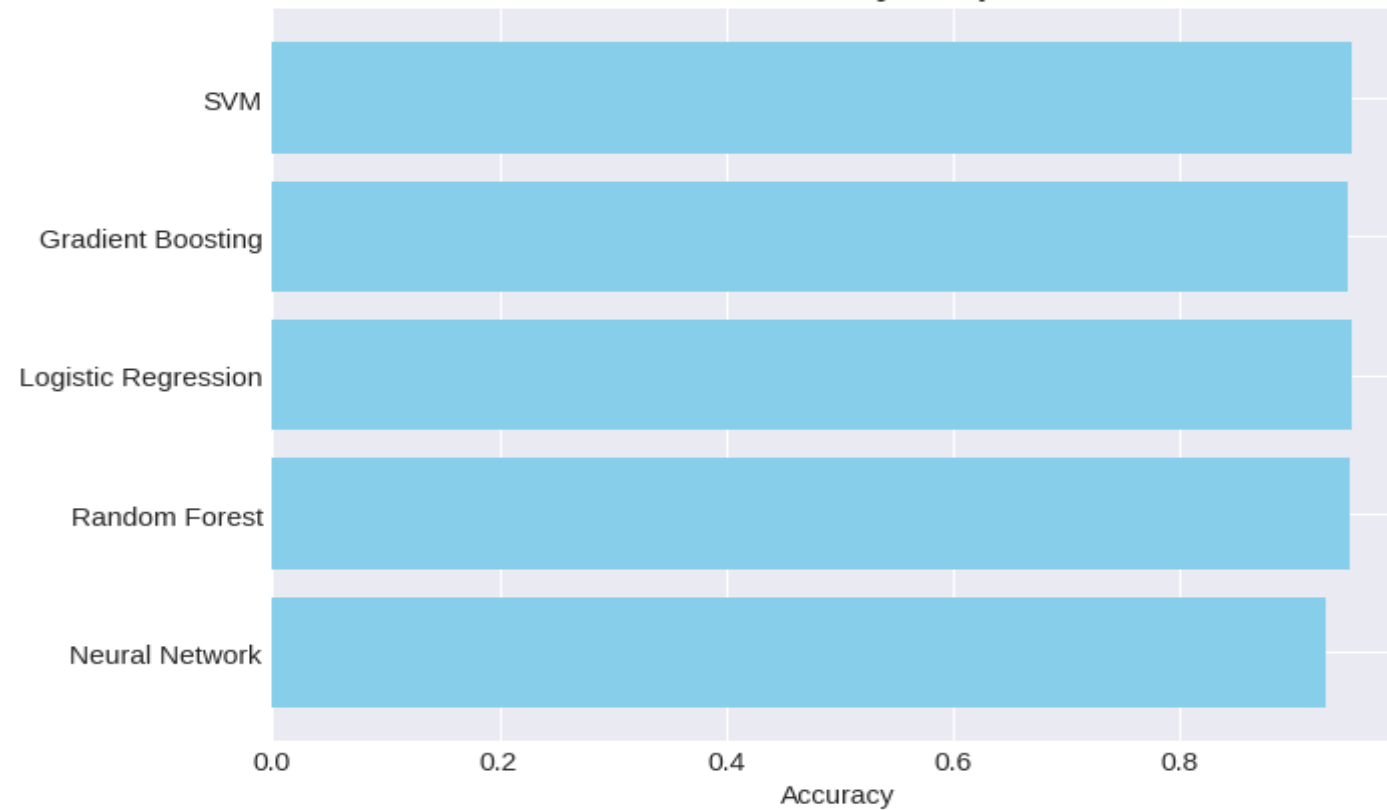


Validation Accuracy

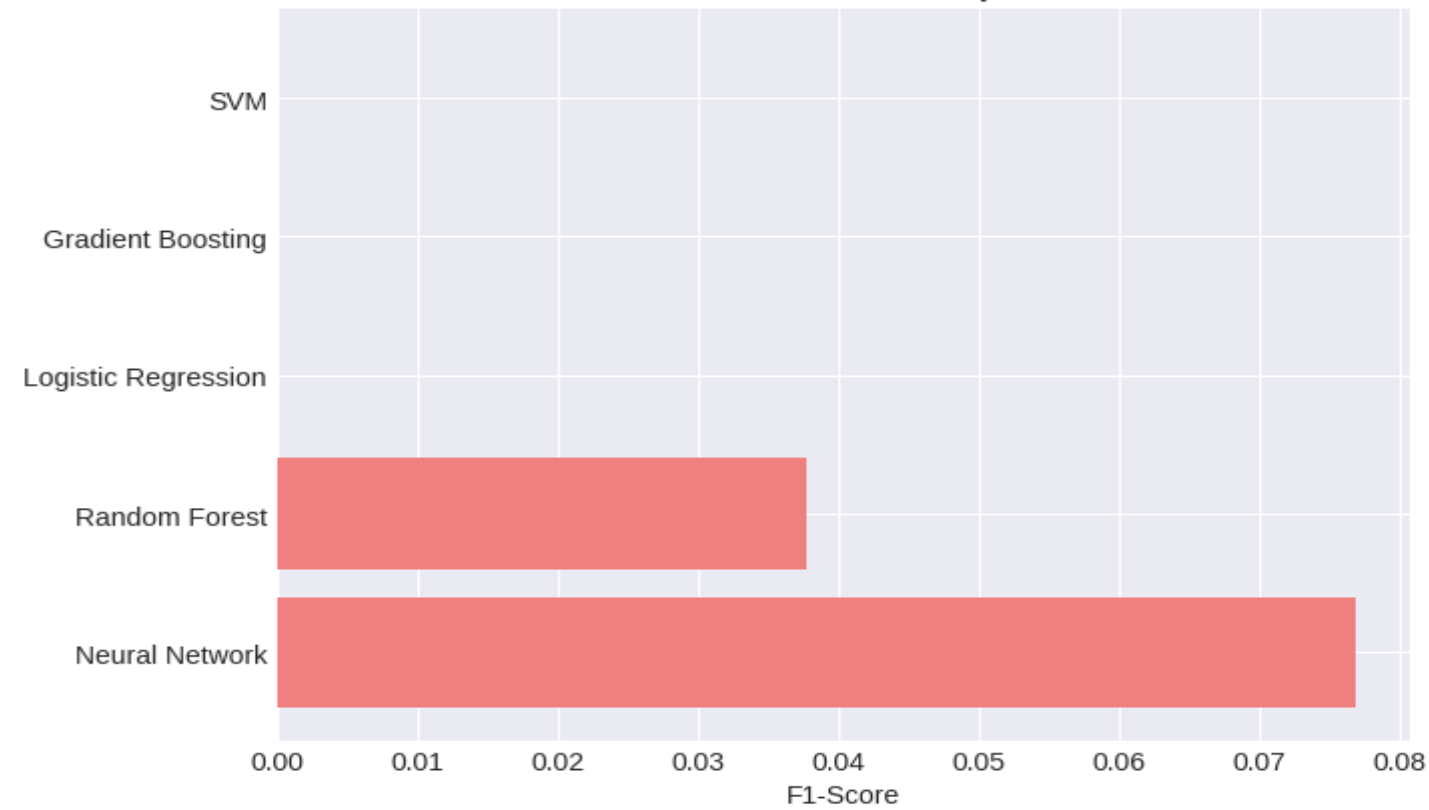
92.95%

Outcome: The model successfully exceeded the 85% accuracy target.

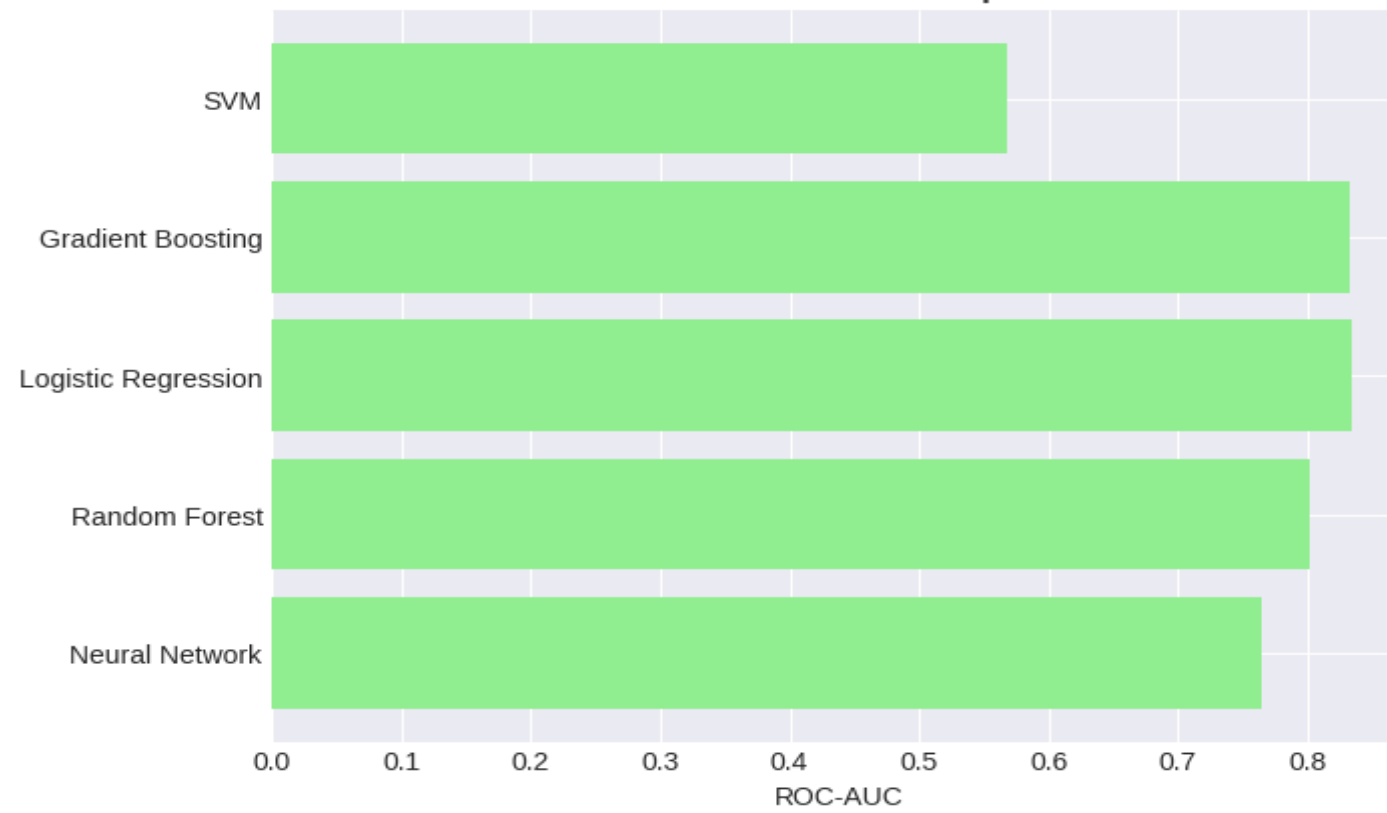
Validation Accuracy Comparison



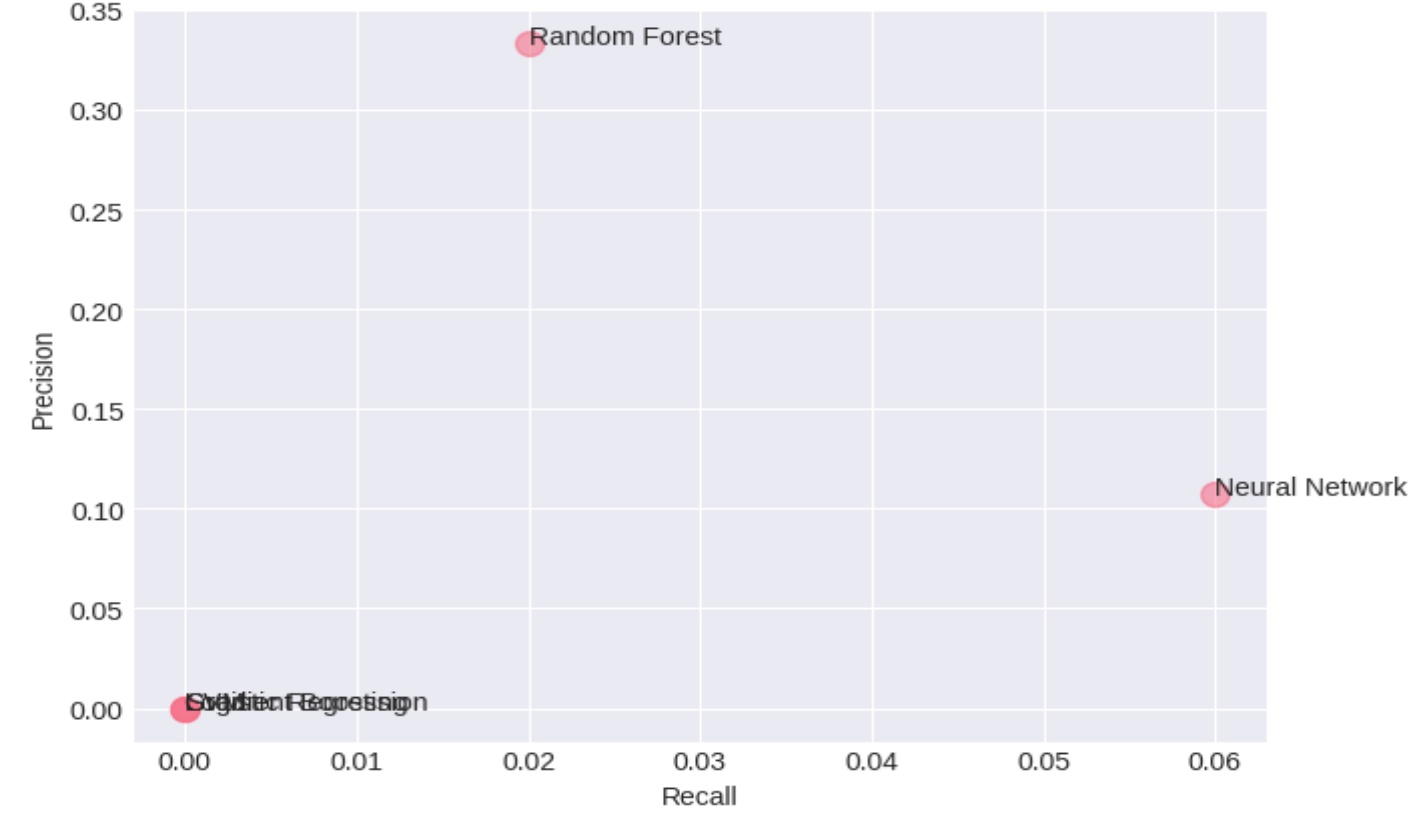
Validation F1-Score Comparison



Validation ROC-AUC Comparison



Precision vs Recall



Deep Dive – Detailed Metrics

F1-Score

Used as a primary indicator to balance precision and recall for the stroke cases.

ROC-AUC

Analyzed the area under the curve to confirm the model's ability to distinguish between stroke and non-stroke patients.

Overfitting & Underfitting Analysis

Metric: Accuracy Difference = 0.0662.

Conclusion:

Since the gap is well below the 0.1 threshold, the model demonstrates "**Good Generalization**," meaning it works reliably on new data.

Testing & Deployment Readiness



Final Testing

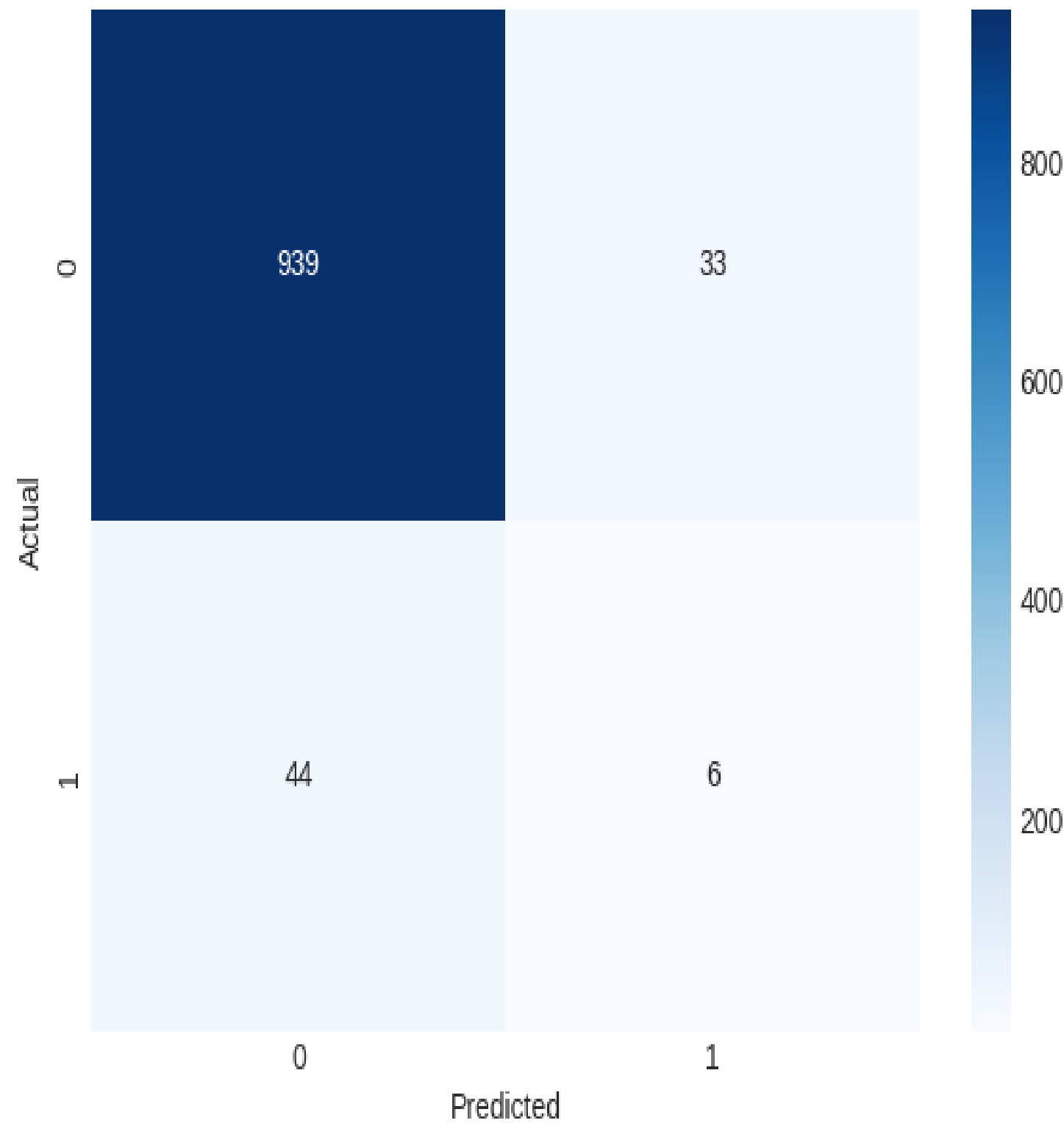
The model was evaluated on the unseen test set to confirm its final capability.



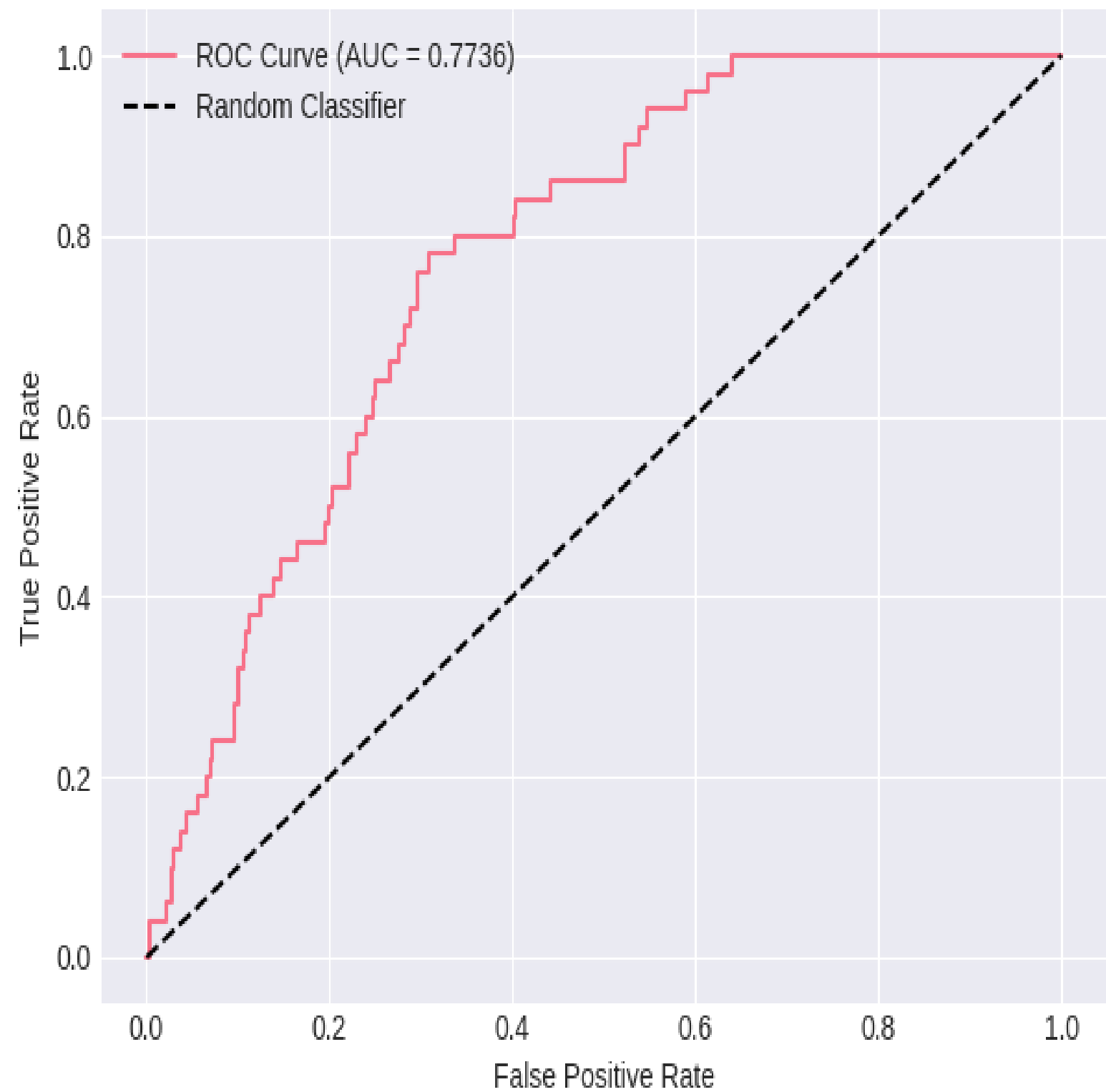
Tools

Saved the final model using joblib or pickle for integration into web apps like Streamlit or Flask.(On progress)

Confusion Matrix

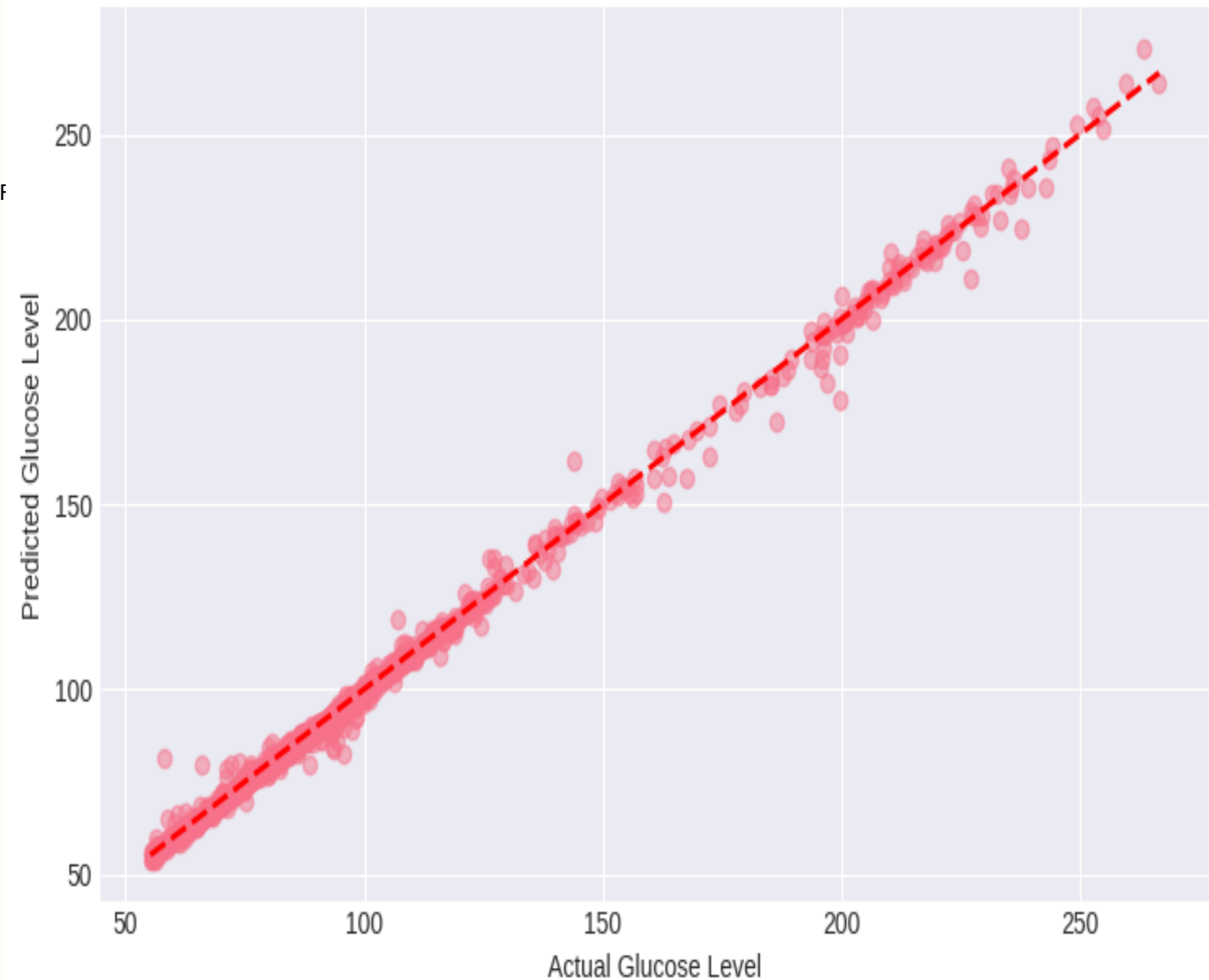


ROC Curve

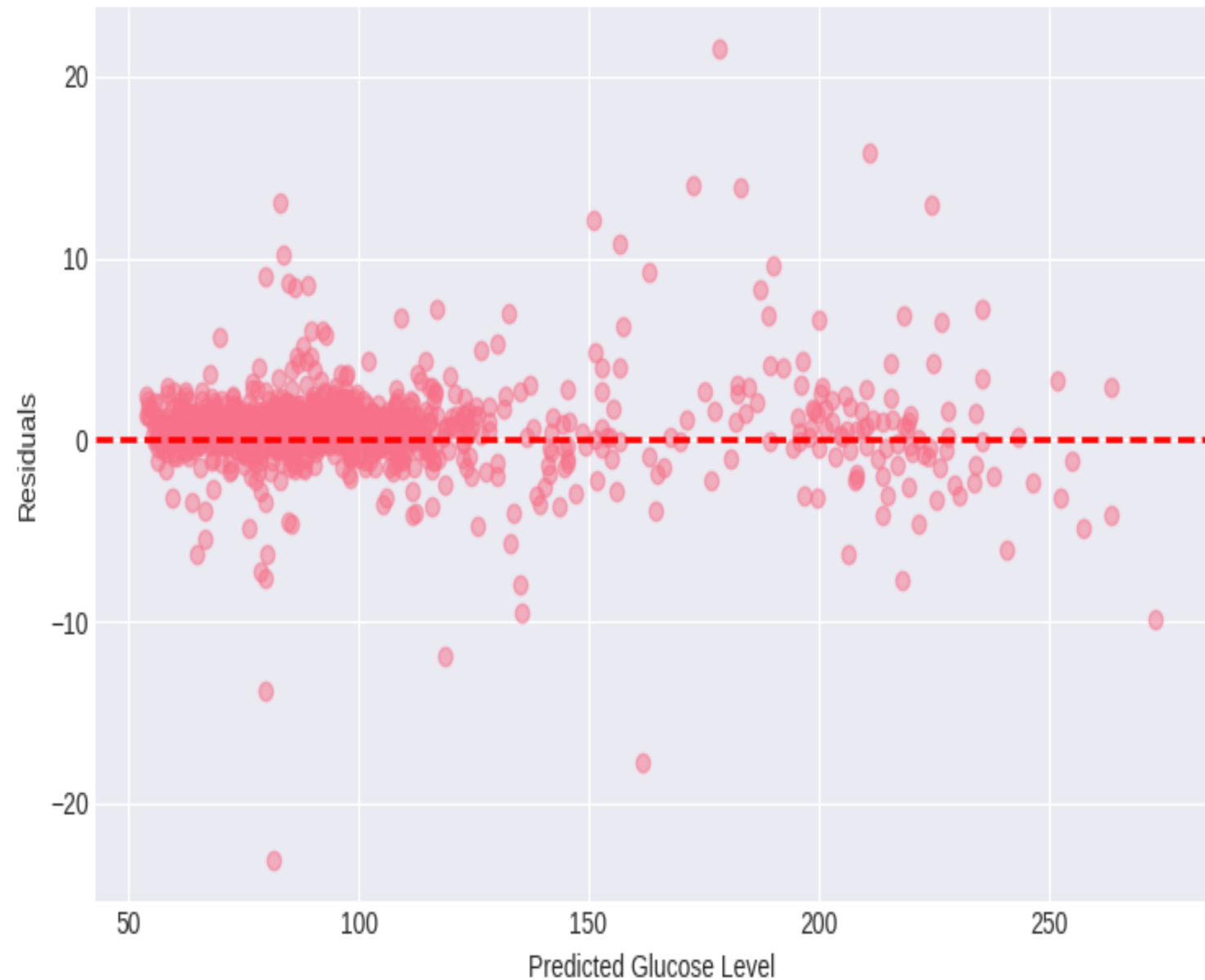


Regression Model - Test Set Results:
RMSE: 2.6226
MAE: 1.5287
R² Score: 0.9964

Predicted vs Actual (R² = 0.9964)



Residuals Plot



Monitoring & Maintenance

Data Drift Monitoring

In a real-world setting, the model would be monitored for "Data Drift" (changes in patient health trends over time).



Continuous Feedback Loop

Continuous data collection would allow for future retraining to improve recall further.

Thank You

Questions & Discussion

Presenter: Natnael Abayneh

