# SANTANDER CUSTOMER TRANSACTION PREDICTION

## TO EVALUATE DIFFERENT MODELS FOR BINARY-CLASSIFICATION

AJEYA KEMPEGOWDA
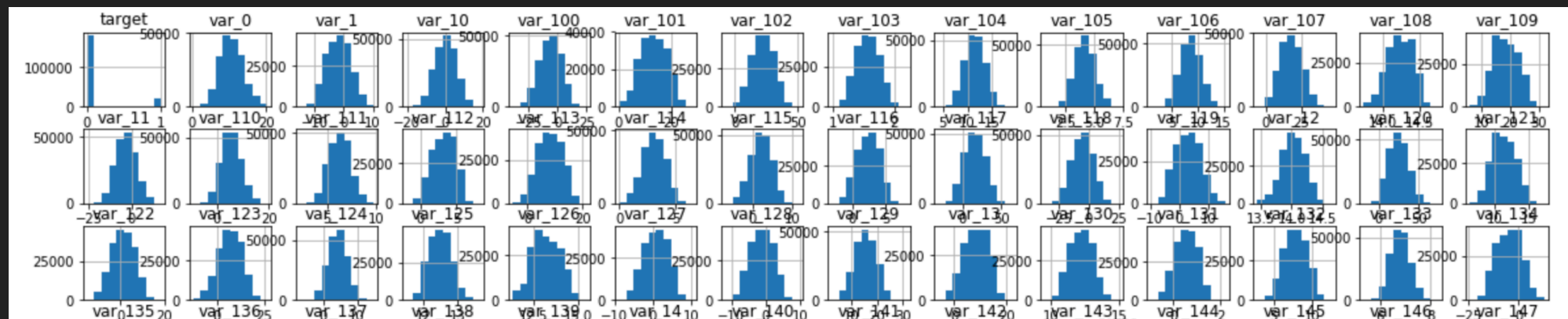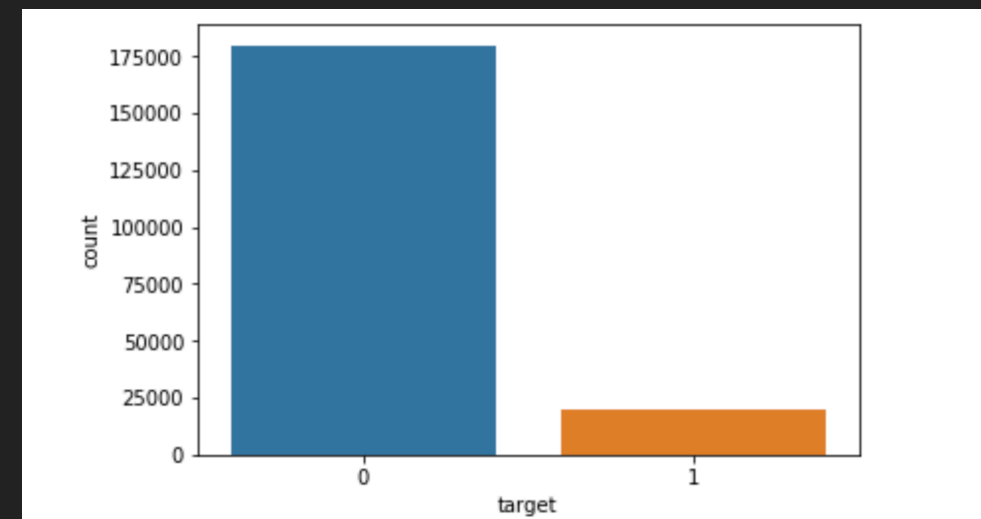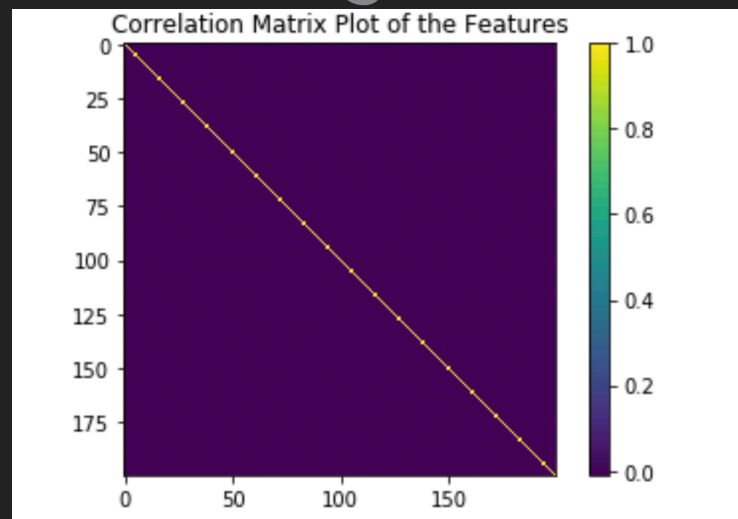DEEPANSHU PARIHAR
HARISH RAMANI
I YEAR MSDS

# PROBLEM STATEMENT

▸ To predict future customer transactions using Santander bank's transaction data.

▸ To build ML models that tackle binary-classification problem.

▸ To explore new family of Neural networks(Neural ODE) and to evaluate the results.

# DATASET

▸ Anonymized real world customer data from Santander bank

▸ Training data available in CSV format.

▸ Training data consists of 200K+ observations and 202 features.

# EXPLORATORY DATA ANALYSIS

▸ Available data is pre-processed - Features are normally distributed

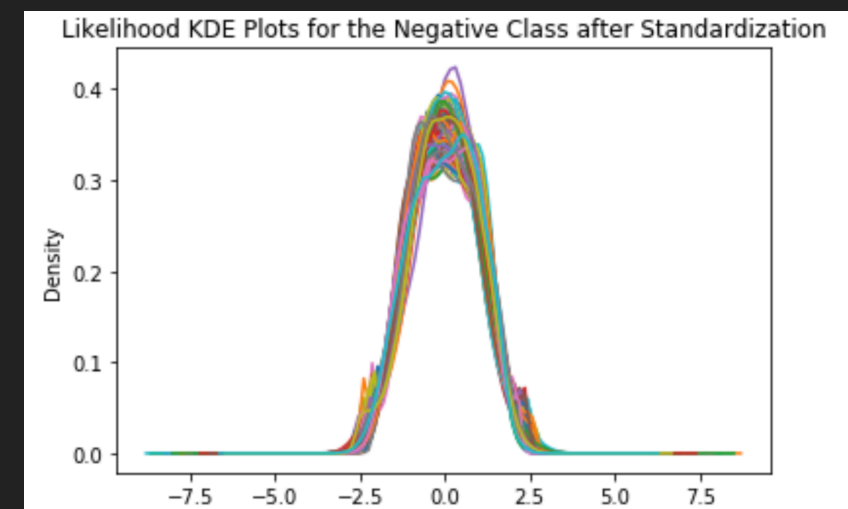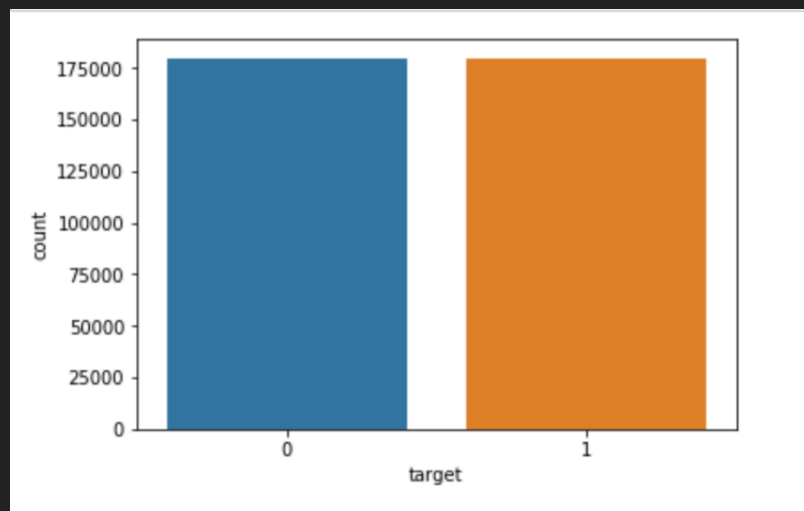▸ Absence of correlation between the features.

▸ Imbalanced target class

# IMPLEMENTATION METHODOLOGY : DATA PRE-PROCESSING

▸ Handling imbalanced data using SMOTE.

▸ Scaled features.



```
Before OverSampling, counts of label '1': 20098
Before OverSampling, counts of label '0': 179902

After OverSampling, the shape of train_X: (359804, 200)
After OverSampling, the shape of train_y: (359804,)

After OverSampling, counts of label '1': 179902
After OverSampling, counts of label '0': 179902
```

# MODEL DEVELOPMENT

▸ **Base line models**

    ▸ Logistic Regression - initial score of 0.630

    ▸ EDA pointed towards Gaussian Naive Bayes - initial score of 0.887

▸ **Neural Networks**

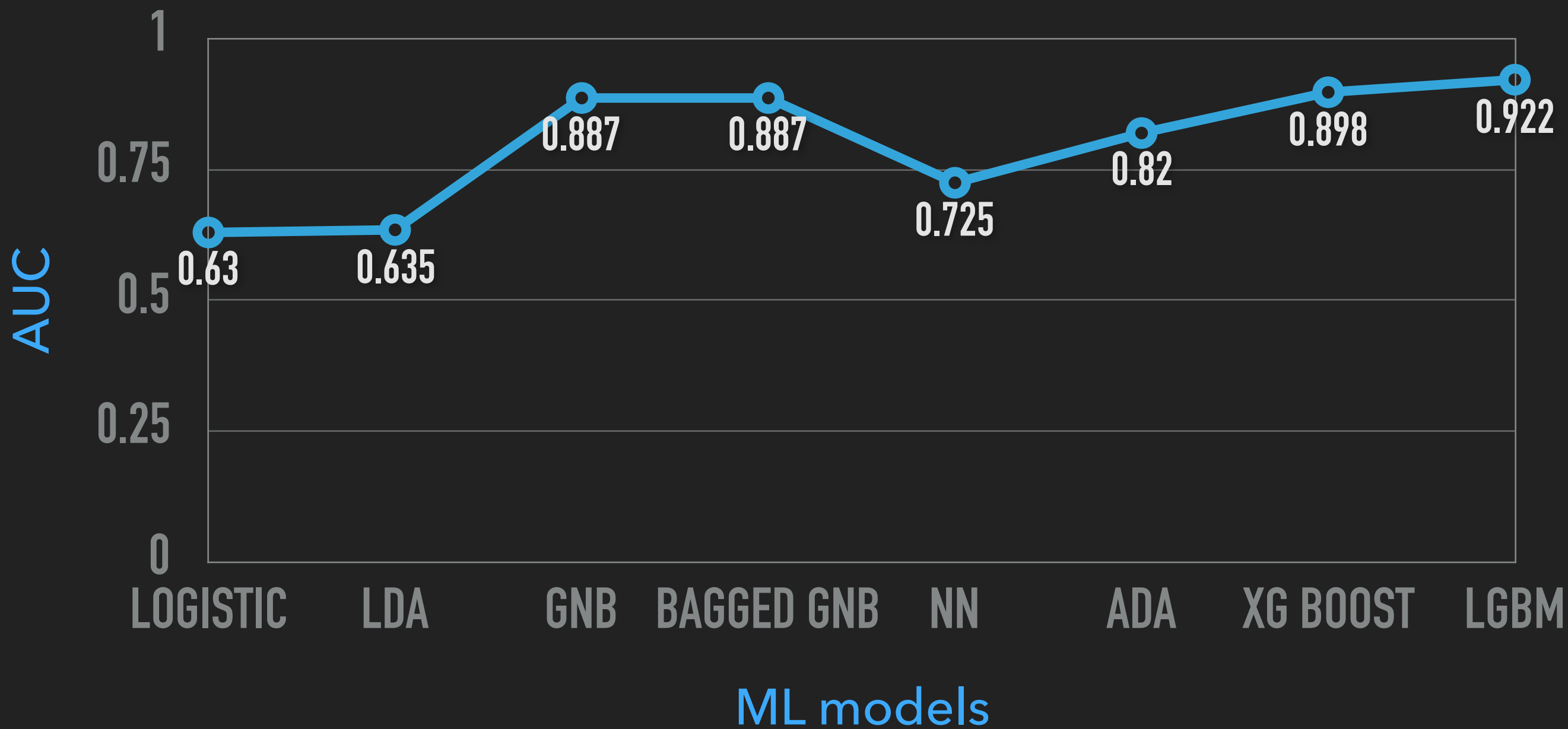    ▸ Feedforward Neural Network **-** Best of 0.725

▸ **Ensemble methods**

    ▸ Gaussian Naive Bayes - - Best of 0.888

▸ **Boosting Methods**

    ▸ XGBoost, LGBM - Best of 0.922 AUC

# MODEL DEVELOPMENT

▸ Feed forward Neural Networks

  ▸ ReLU activation

  ▸ Learning rate = 1e-4

  ▸ Iterations = 1000

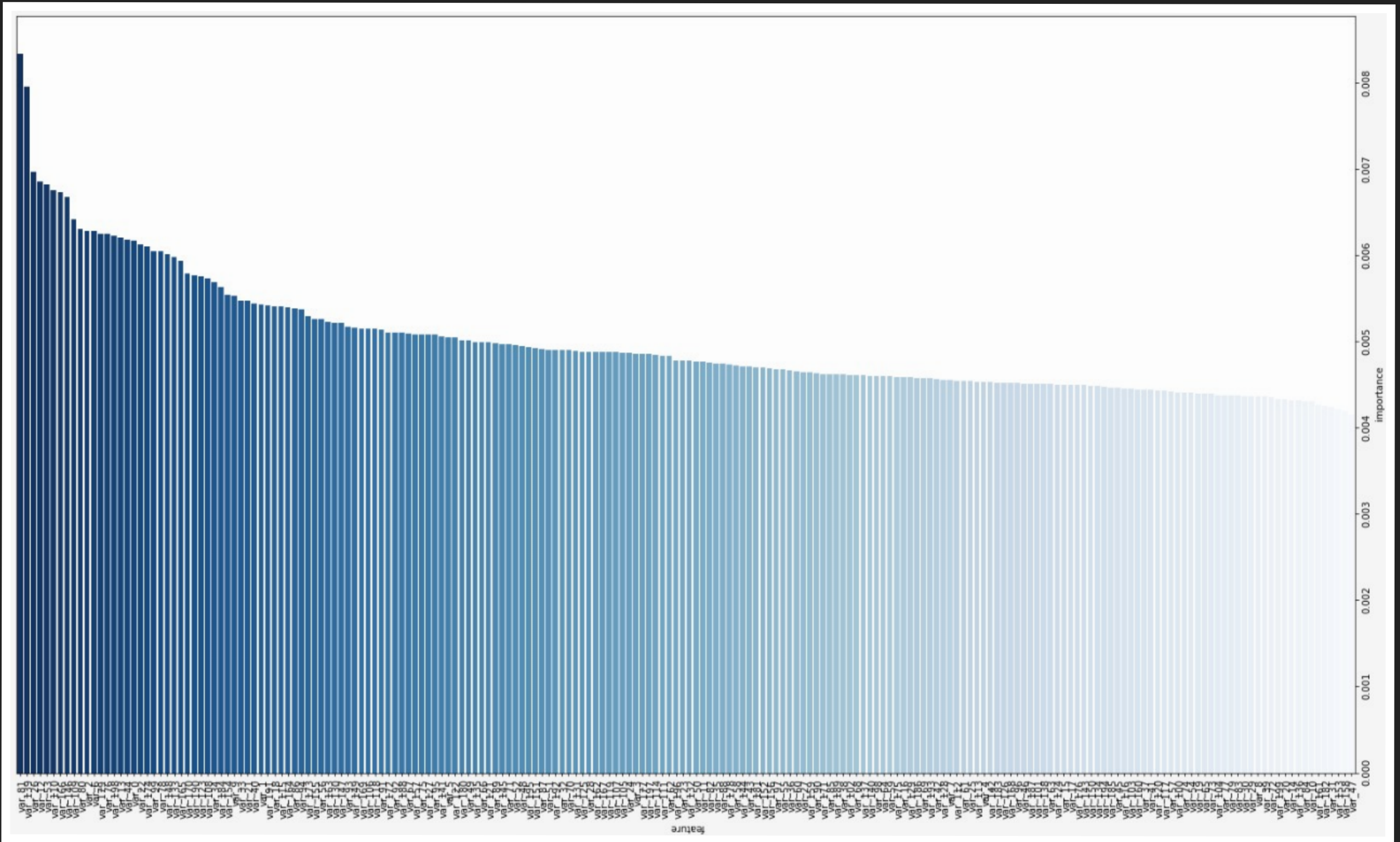  ▸ Hidden layers = 3(x 13 neurons each)

▸ Boosting Methods

  ▸ XGBoost

    ▸ Learning rate = 0.01

    ▸ Max depth = 2

    ▸ Learning objective = binary logistic

  ▸ LGBM

    ▸ Boosting type = Gauss

    ▸ Learning rate = 0.05

    ▸ Max depth = 5

▸ Parameter tuning using Grid search

# FEATURE IMPORTANCE

# KAGGLE SUBMISSIONS AND SCORE EVALUATIONS

# NEURAL ODE – INTRODUCTION

▸ Parameterizes the derivative of hidden state layers.

▸ Provides a continuous depth model.

▸ Properties we understood

  ▸ Memory efficiency: Constant memory cost wrt depth

  ▸ Adaptive computation: Adapt error levels for accuracy

  ▸ Continuous time series model: For time series model (unlike RNN's)

# NEURAL ODE – IMPLEMENTATION

▸ 'odeint' interface to solve the initial value problem(ODE+initial state)

▸ odeint(ODE solver ) tries to find the trajectory satisfying the ODE that passes through the initial conditions.

▸ ODE solver can be tweaked to acquire fixed steps(Euler) or to adaptive(Runge Kutta)

▸ Back-propagation is done using 'ode_adjoint' that solves adjoint ODE in O(1) space complexity.