# Project Report

# Predicting Customer Churn

Shayan Lodi - 27002
Faariha Atiq – 26993
Uzair Arif – 27184
Kanza Rehan - 27182

## Business Problem

One of the primary objectives for any business is to mitigate customer churn, as it directly impacts revenue and incurs additional costs in acquiring new customers. As a data scientist at a bank, the task at hand is to leverage historical data to forecast whether a customer will churn in the upcoming six months. This predictive model will enable the bank to strategize proactive engagement with customers, ensuring timely interventions to retain them effectively.

## About the data

### Variables

Each row in the dataset is uniquely identified by an ID. The age of each customer is recorded, along with their gender, which can be either male or female. The yearly income of the customer is also recorded. Additionally, the average quarterly balance of the customer's account and the length of time they have been associated with the bank, referred to as "Vintage," are included. The dataset indicates whether the customer has conducted any transactions in the past three months. Furthermore, it records the number of product holdings each customer has with the bank and whether they possess a credit card. The credit category of each customer, determined by their credit score, is also noted. Finally, the dataset specifies whether each customer is predicted to churn within the next six months.

## Normalisation

Normalization is a technique used to scale numerical features to a similar range. This ensures that features with different scales do not disproportionately influence the model's training process.

## Discretization

Discretization involves converting continuous numerical features into categorical features by dividing them into intervals or bins. This can help simplify the model and make it more interpretable, especially when dealing with algorithms that require categorical input.
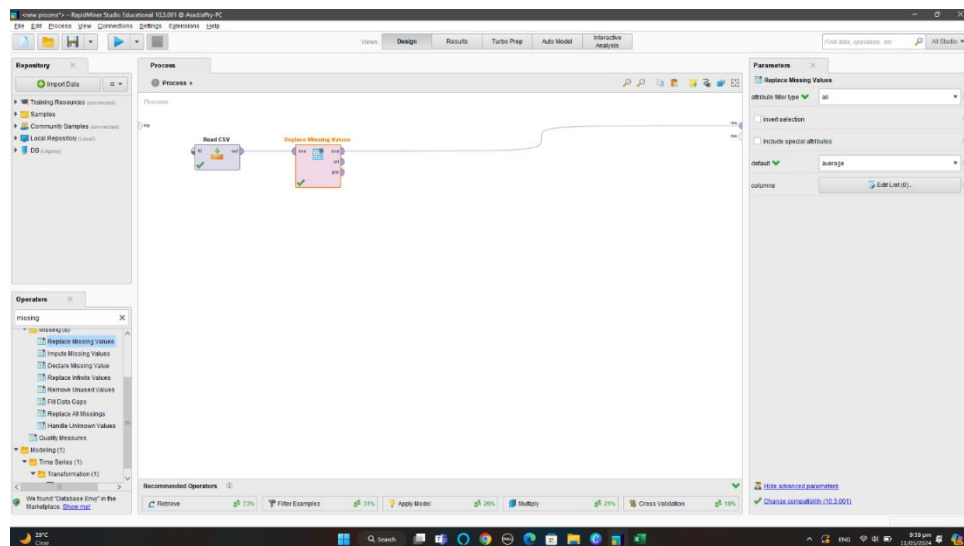
## Feature Encoding

Feature encoding is the process of converting categorical features into a numerical representation that machine learning algorithms can understand. Feature encoding ensures that categorical data can be effectively utilized in the modeling process.

# Cleaning the data

We start by first cleaning our data. Cleaning data before applying any data science techniques is essential to ensure accuracy, completeness, consistency, relevance, and quality of the dataset. It involves removing any errors, handling missing values, standardizing formats, and eliminating irrelevant variables. In essence, data cleaning is a critical preparatory step that enhances the integrity and usefulness of the data for subsequent data science tasks.

To handle missing values, we used the Replace Missing Values operator and applied average value option for the missing values.

The result is as follows:



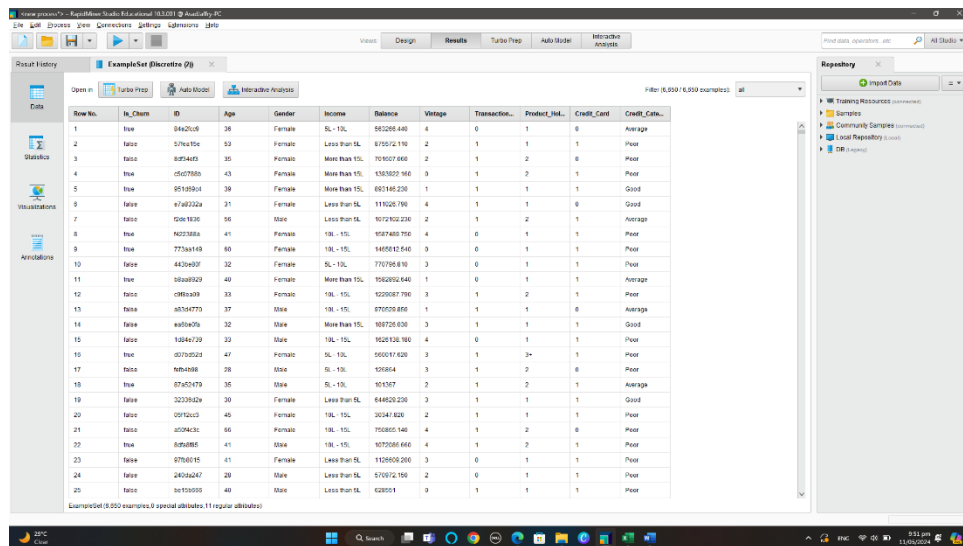The target variable 'Is_Churn' is already normalized so we don't need to apply normalization.

# Logistic Regression

Logistic regression is a statistical method used for binary classification tasks, predicting the probability of an outcome based on predictor variables. It's chosen for its interpretability, probabilistic outputs, computational efficiency, and flexibility in handling nonlinear relationships. Logistic regression provides valuable insights into binary outcomes and is widely utilized in data analysis and predictive modeling.

To prepare the data for logistic regression, we use the Discretize by User Specification operator. This involves categorizing the data into distinct classes based on predefined criteria. In this case, we assign the label 'false' to values of 0 and 'true' to values of 1. This process simplifies the data, making it suitable for logistic regression analysis.
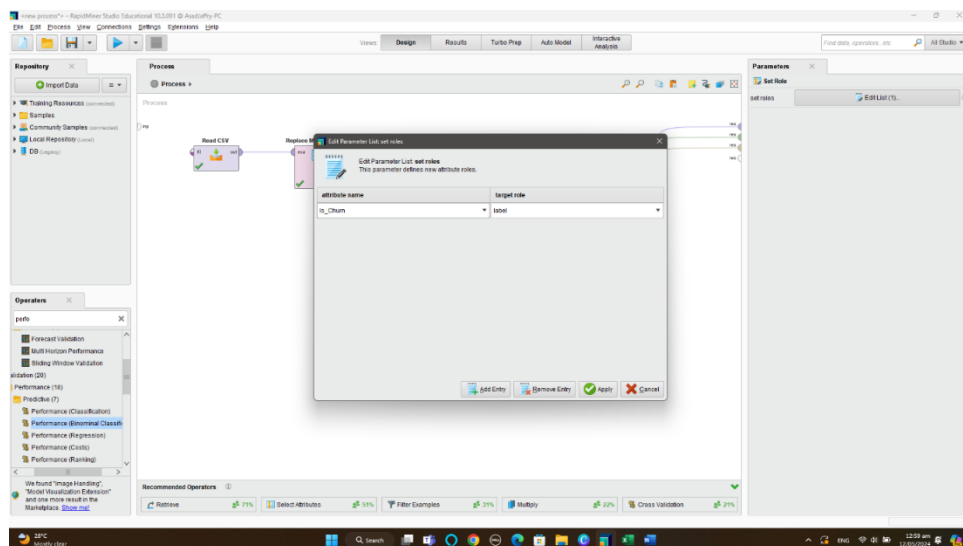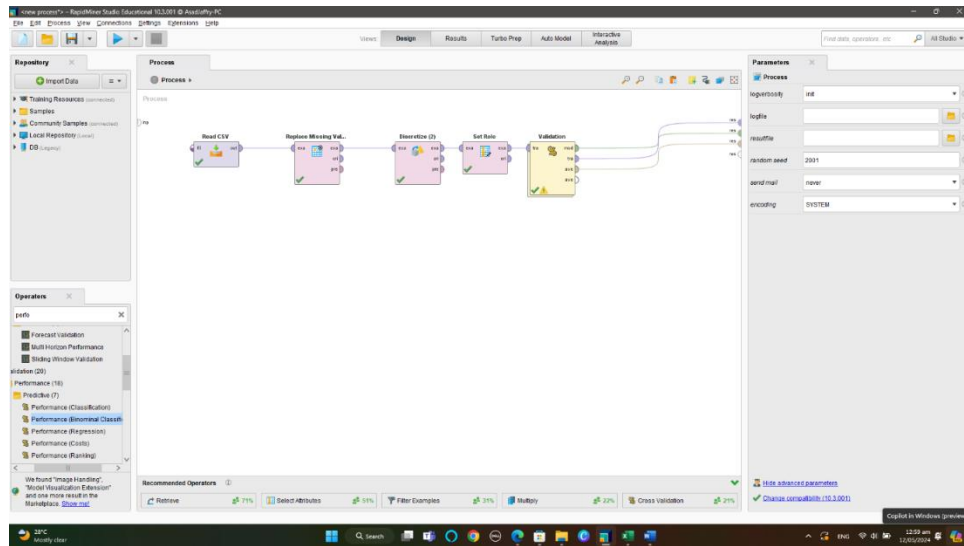
The result is as follows:



Next, we use the Set Role operator to set 'Is_Churn' as 'label'.



We then use Split Validation to split our data into a 70:30 ratio using stratified sampling. Stratified sampling is a sampling technique used to ensure that the proportions of different classes or categories within the dataset are maintained in the sample. This is particularly useful when dealing with imbalanced datasets, where one class may be significantly more prevalent than others.

Double clicking the validation operator, we go to its subprocess. Here we add the Logistic Regression operator, Apply Model and Performance (Binomial Classification) operators. For performance, we choose accuracy, recall, precision, and AUC as shown below.

The resultant Confusion Matrix and AUC curve is this:





Accuracy comes out to be 23.07%, AUC = 0.5 (ideal), precision is 23.07%, and recall is 100%.

# Decision Tree

A decision tree is a machine learning algorithm used for classification and regression tasks, organizing data into a tree-like structure based on pre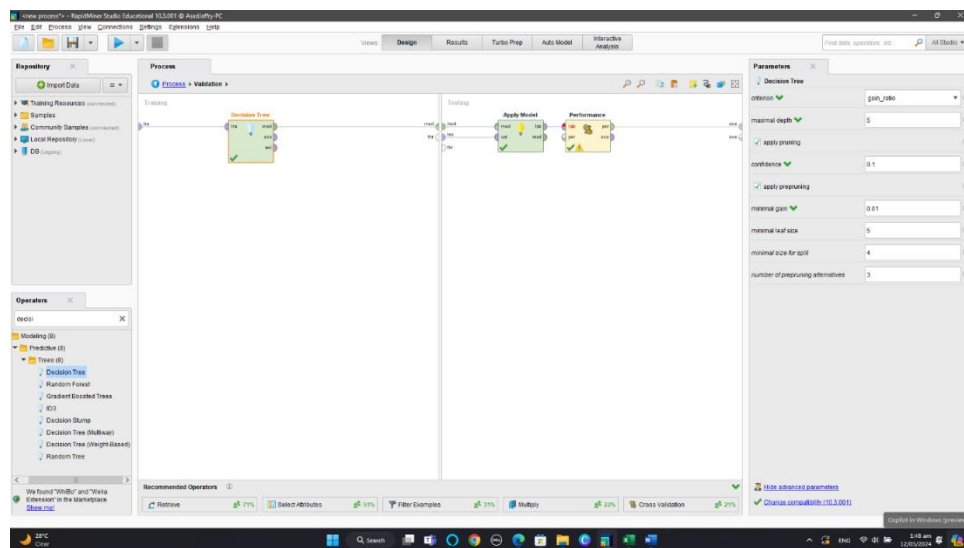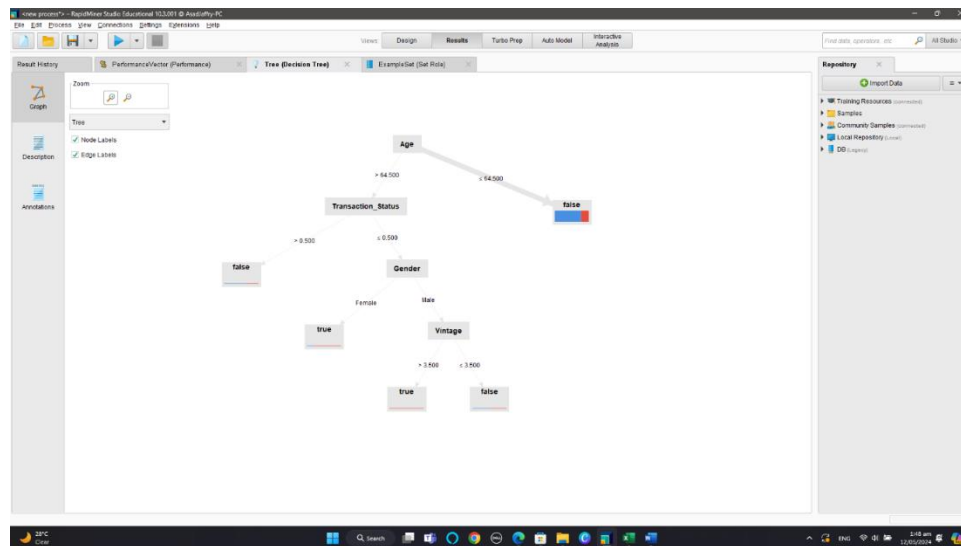dictor variables. We employ decision trees on our data for their interpretability, ability to capture nonlinear relationships, feature importance ranking, versatility, and scalability, providing actionable insights for decision-making tasks.

Using the Decision Tree Operator, with parameters 5,5, and 4 for maximal depth, minimal leaf size and minimal size for split respectively. The maximum depth parameter specifies the maximum number of levels allowed in the decision tree. The minimal leaf size parameter determines the minimum number of samples required to consider creating a leaf node. The minimal size for split parameter specifies the minimum number of samples required to split a node further.

We get the following results.

1. Decision tree



2. Confusion Matrix



Accuracy has now increased to 76.93% compared to logistic regression. The increase in accuracy from logistic regression to decision tree suggests that the decision tree model might be better capturing the underlying patterns and relationships in the data compared to logistic regression. This improvement could be attributed to the decision tree's ability to handle nonlinear relationships and interactions between variables more effectively.

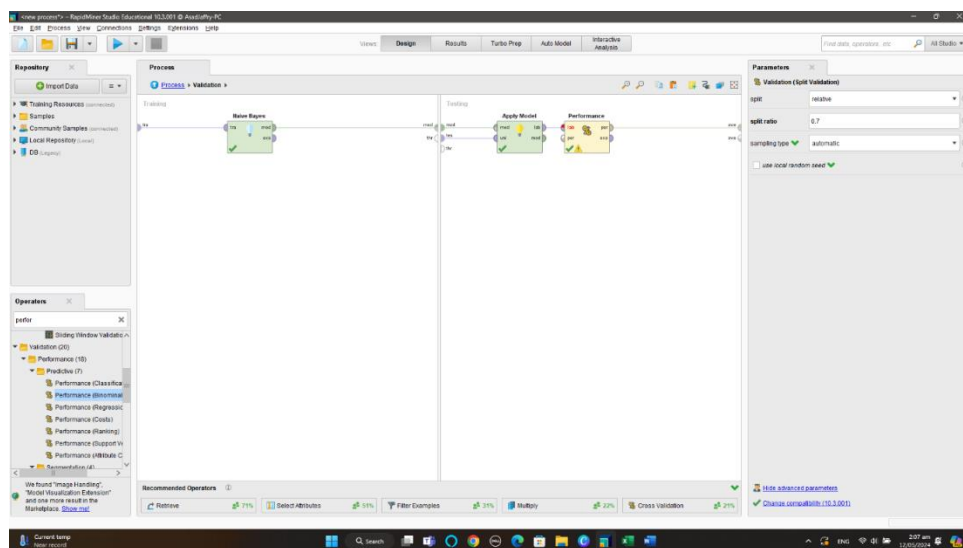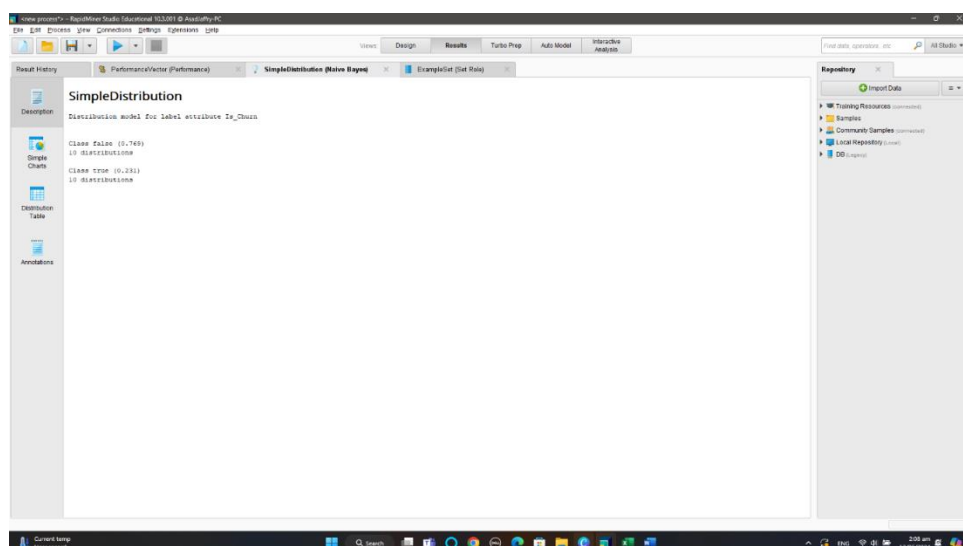3. AUC

AUC equals 0.5, which is ideal.

# Naïve Bayes

Naive Bayes is a probabilistic machine learning algorithm utilized primarily for classification tasks. It operates on the assumption of feature independence given the class variable, making it computationally efficient and easy to interpret. We apply Naive Bayes to our data due to its simplicity, efficiency, and versatility in handling both numerical and categorical data. Its scalability and suitability for large datasets make it particularly useful for real-time applications and big data scenarios. Additionally, Naive Bayes serves as a baseline model for comparing the performance of more complex algorithms. Overall, Naive Bayes is a valuable tool in data analysis and predictive modeling, offering a balance of simplicity and effectiveness.
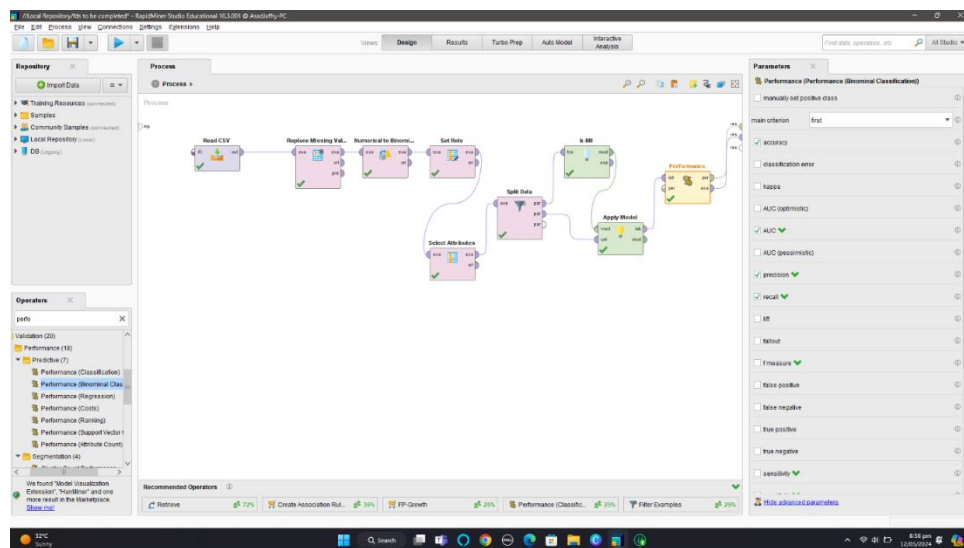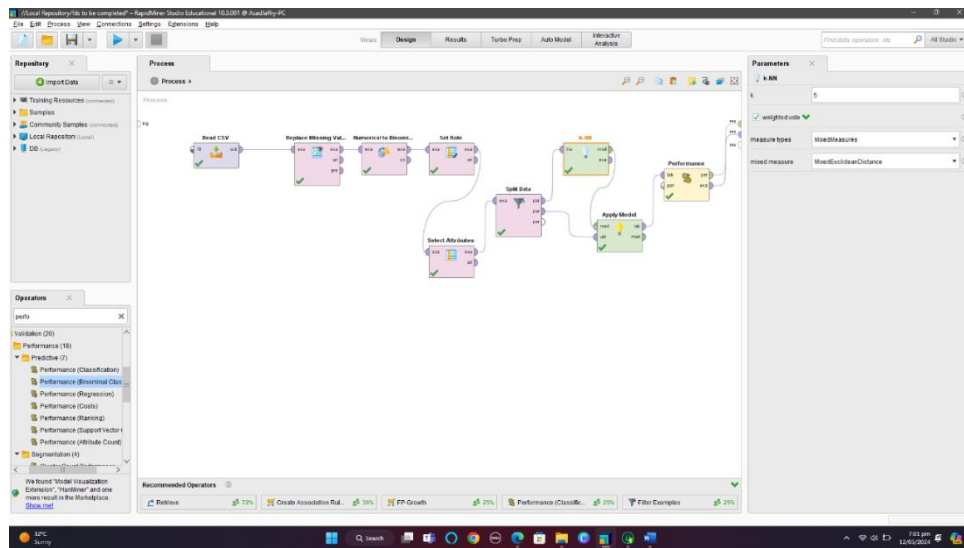


Result:

AUC for Naïve Bayes is 0.635, higher than the logistic regression and decision tree. The higher AUC indicates that Naive Bayes is better at distinguishing between the positive and negative classes in the dataset. This suggests that Naive Bayes may have learned the underlying patterns and relationships in the data more effectively, resulting in improved predictive performance.

# KNN

K-Nearest Neighbors (KNN) is a versatile supervised machine learning algorithm utilized for classification and regression tasks. Its simplicity and ease of implementation make it a popular choice, particularly for initial exploration of datasets. KNN operates by assigning a new data point to the majority class (for classification) or computing the average of its nearest neighbors (for regression) based on a specified number, K, of closest data points in the feature space. We perform KNN on our data for its flexibility in handling various types of predictive modeling problems, as well as its ability to capture complex relationships between variables without making assumptions about the underlying data distribution. Additionally, the interpretability of KNN results, coupled with its lack of a training phase, makes it an appealing option for quick insights and decision-making tasks in data analysis.



Parameters for KNN; k=5. What this means is that the algorithm considers the five nearest neighbors when making predictions for a new data point. These neighbors are determined based on Euclidean distance, and the majority class among the nearest neighbors is used as the predicted class for classification tasks.

Result:





Accuracy for KNN is 72.17%, higher than the previous three machine learning models.

AUC is 0.508, and the precision is 24.04%.

# SVM

Support Vector Machine (SVM) is a versatile supervised machine learning algorithm used for classification, regression, and outlier detection tasks. SVM identifies the hyperplane that best separates data points into different classes, aiming to maximize the margin between them. We 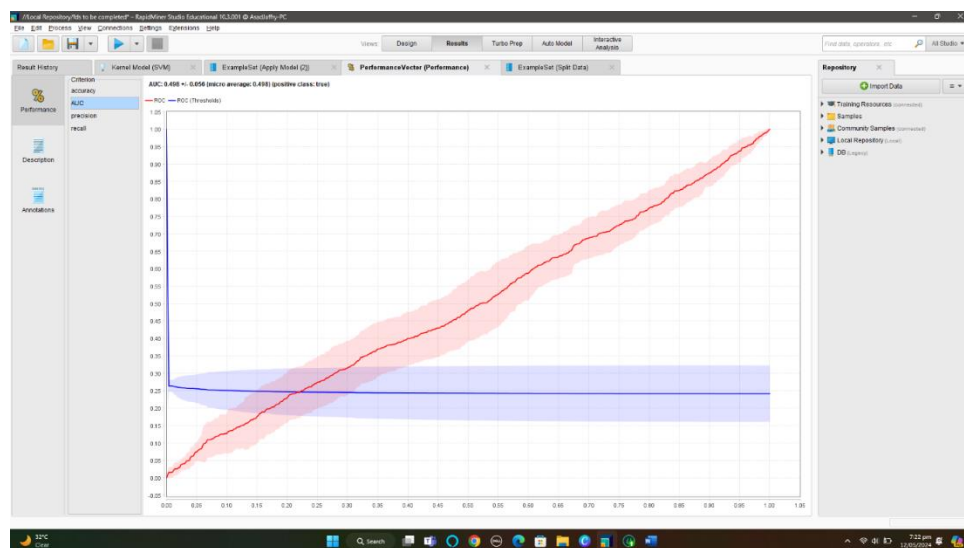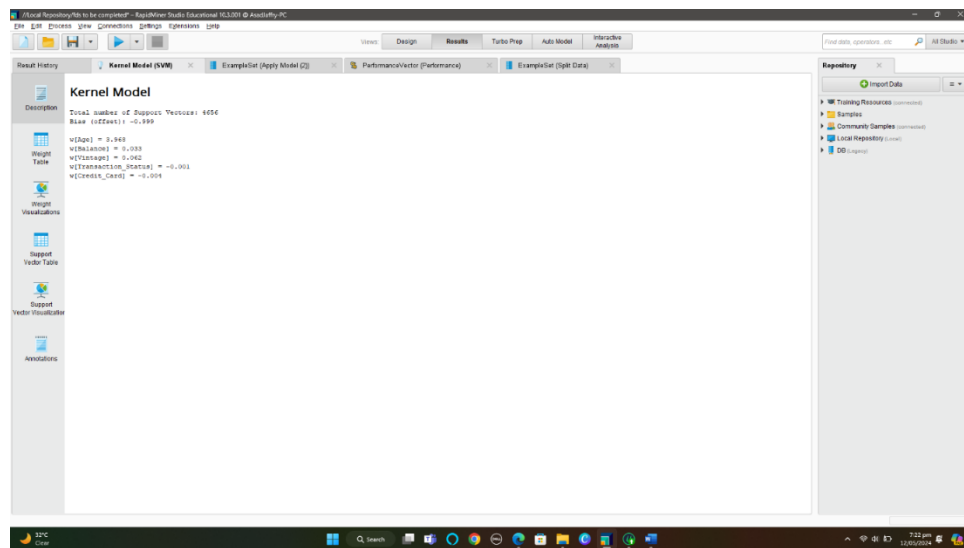utilize SVM on our data due to its versatility in handling both linear and nonlinear classification tasks through various kernel functions. SVM is robust to overfitting, particularly in high-dimensional spaces, making it suitable for datasets with many features. Its effectiveness in high-dimensional spaces, along with its focus on maximizing the margin between classes and robustness to outliers, makes SVM a valuable tool for various data analysis and predictive modeling tasks.
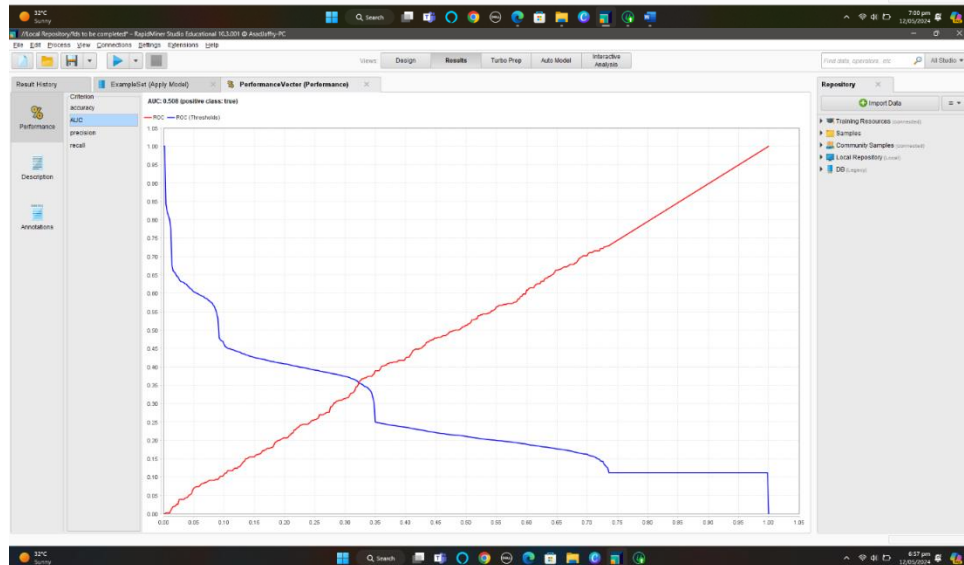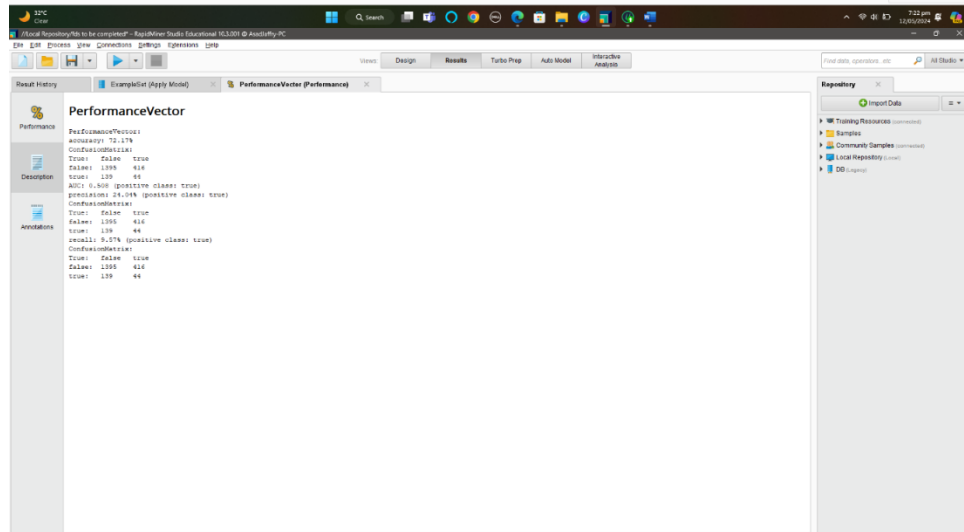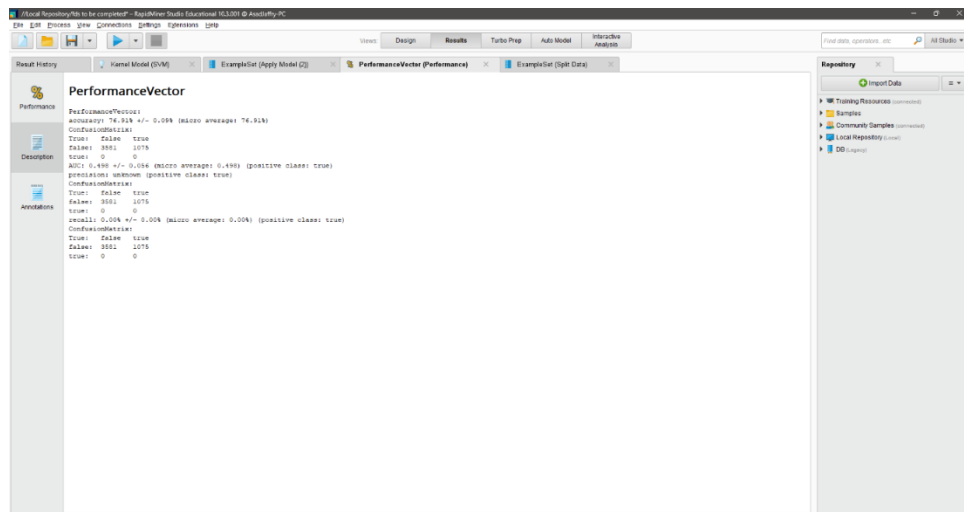
Kernel type is set to polynomial so to identify any complex, non-linear relationships between input features and the target variable. It allows the machine learning model to project the data into higher-dimensional spaces, where it becomes easier to find a hyperplane that can effectively separate the classes.

$C$ is a hyperparameter known as the regularization parameter. It controls the trade-off between maximizing the margin and minimizing the classification error. Setting $C$ to zero in SVM effectively removes the regularization term from the objective function. This means that the SVM algorithm will focus solely on maximizing the margin, without penalizing misclassifications.
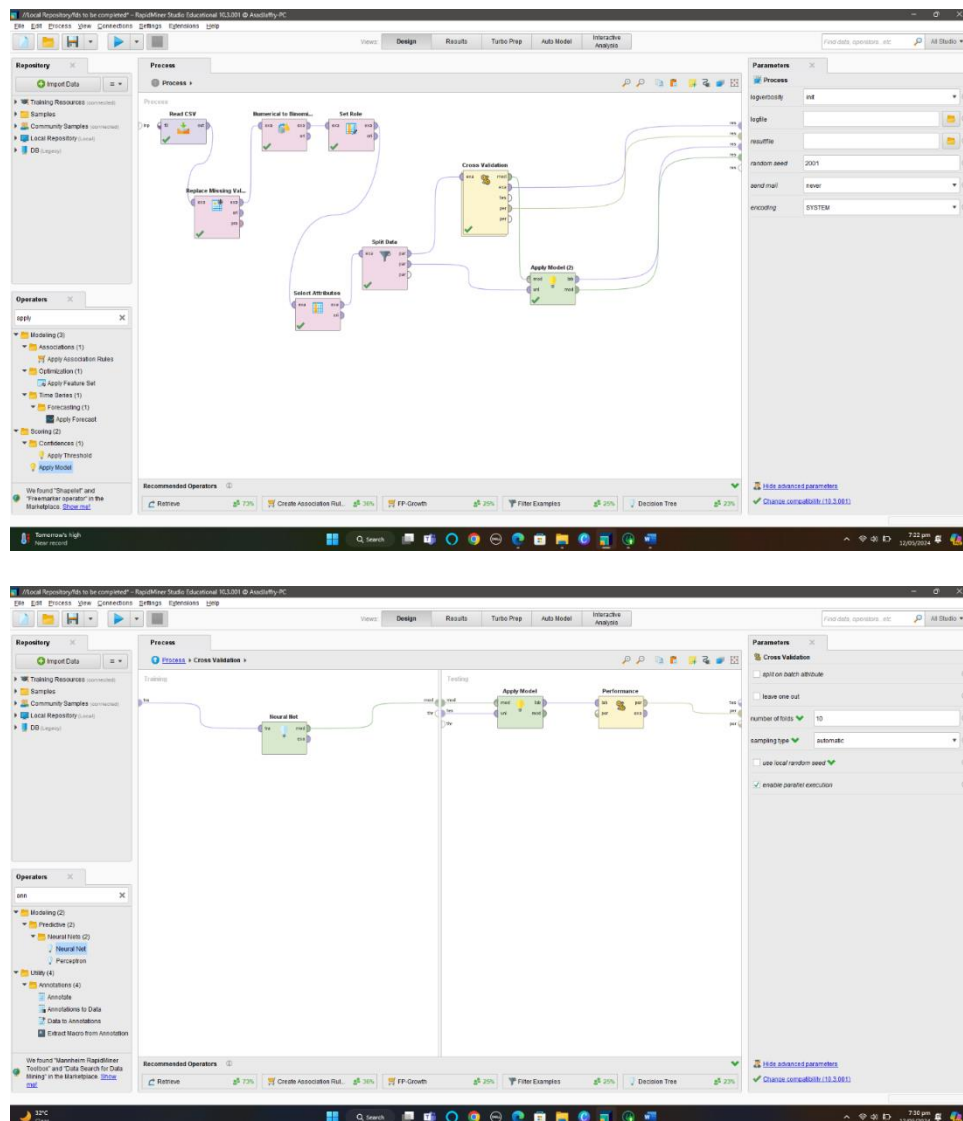




AUC for the kernel model comes out to be 0.498.

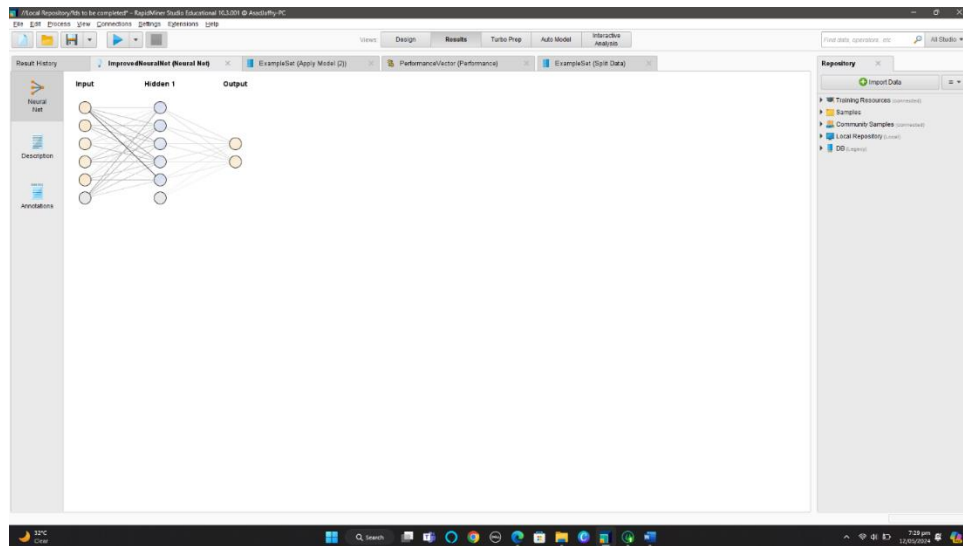Accuracy for performance vector is 76.91%. AUC is 0.498

# ANN

Artificial Neural Networks (ANNs) are sophisticated machine learning models inspired by the structure and function of the human brain. Comprising interconnected nodes organized into layers, ANNs learn complex patterns and relationships in data through training, adjusting weights and biases based on input-output pairs. We employ ANNs on our data due to their capability to capture nonlinear relationships, automate feature learning, and achieve state-of-the-art performance across various tasks such as classification, regression, and unsupervised learning.
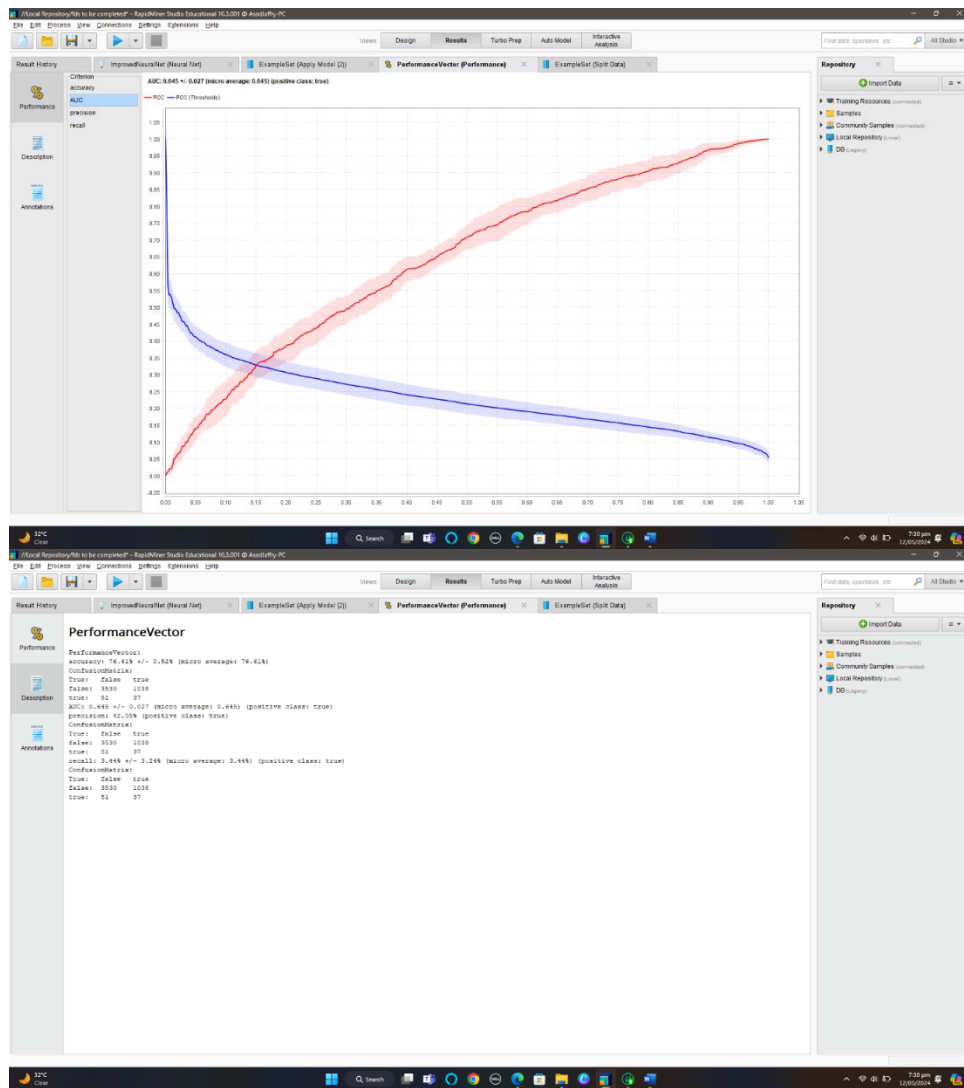
Setting the sampling type to "automatic" for Artificial Neural Networks (ANNs) indicates that the algorithm will automatically handle the sampling strategy based on the dataset characteristics and the specific implementation of the ANN algorithm being used. This simplifies model selection and potentially improves performance.

Result:

Accuracy for ANN model is 76.61%. AUC comes out to be 0.645 and precision is 42.05%.

# Comparison between models

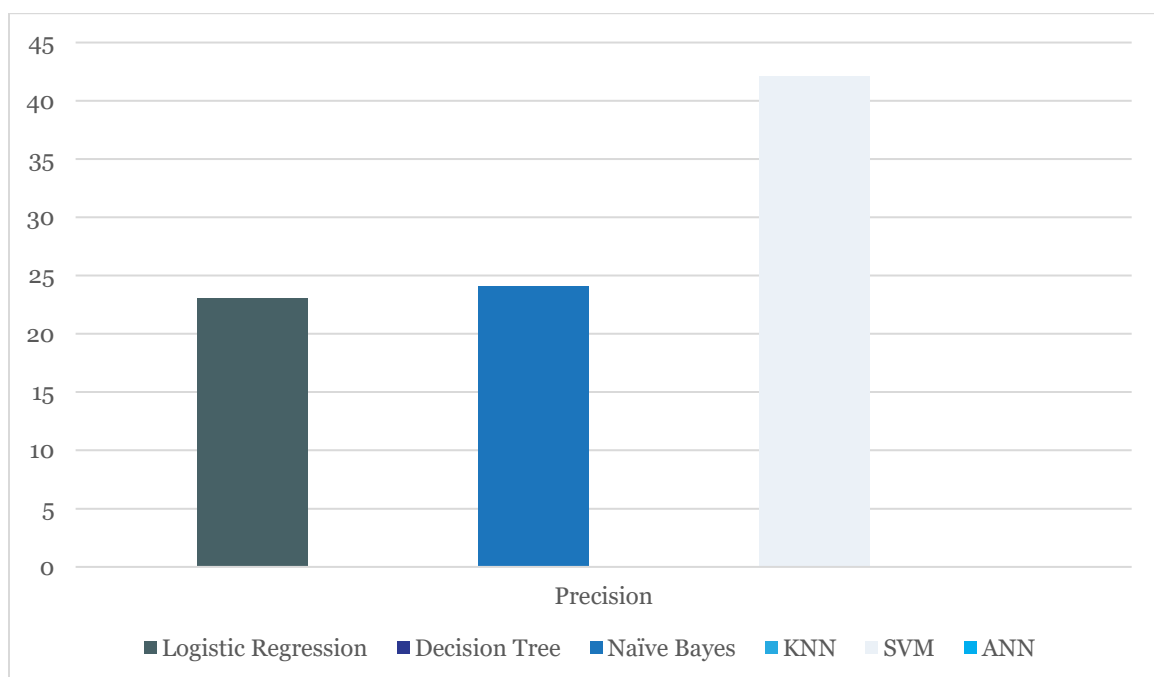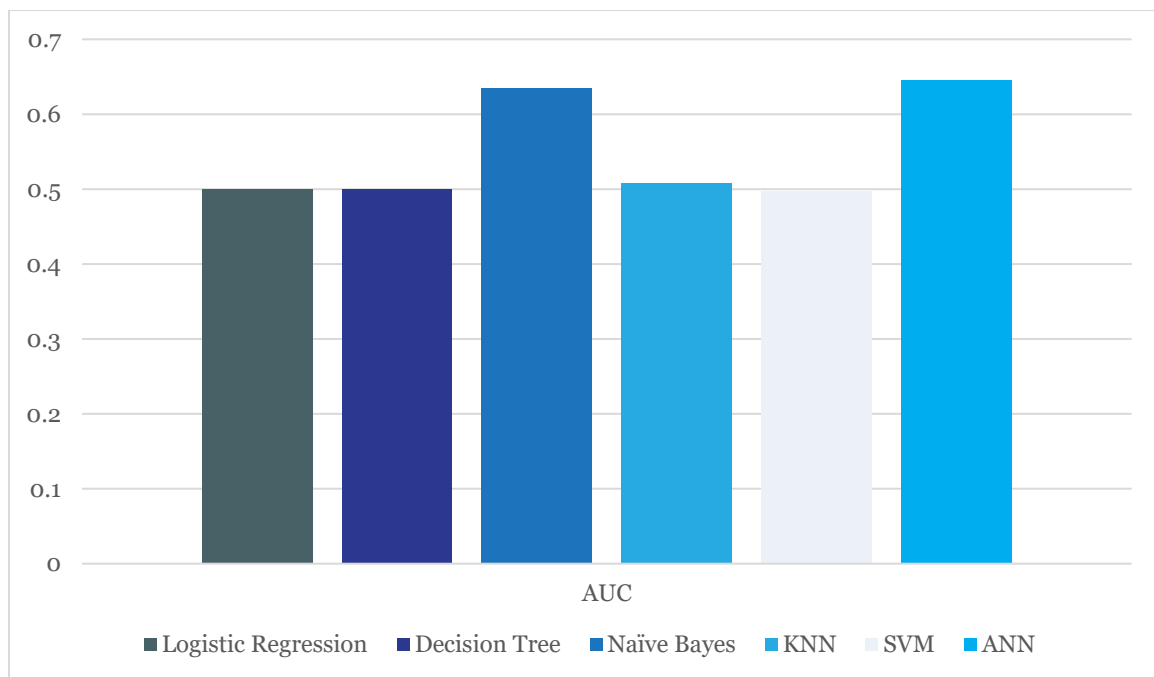Comparing the performance of machine learning models is essential for evaluating their effectiveness in classification tasks. Accuracy measures the overall correctness of predictions, precision quantifies the proportion of correctly predicted positive cases among all predicted positive cases, and recall assesses the proportion of correctly predicted positive cases among all actual positive cases. By visualizing these metrics using bar charts, we can gain insights into how different models perform across multiple evaluation criteria. This comparison allows us to identify the strengths and weaknesses of each model and make informed decisions about which model is best suited for the specific task at hand.

The results show that Decision Tree, SVM, and ANN outperform Logistic Regression, Naïve Bayes, and KNN in terms of accuracy. However, it's essential to note that Decision Tree has a recall of 0%, indicating potential issues with correctly identifying positive cases. Naïve Bayes demonstrates the highest AUC, suggesting its better performance in distinguishing between positive and negative cases. Despite having a relatively low recall, ANN exhibits the highest precision, indicating its capability to correctly classify positive cases.