



BA 706 Project

# DATASET: RBC CHURN DATASET

Grace Adaji

Chukwudera Stephen Chukelu

Alicia Waters

Chidiebere Odurukwe

## Table of Contents

|   |           |
|---|-----------|
| <b>1. Executive Summary .....</b>             | <b>1</b>  |
| <b>2. Objective .....</b>                     | <b>3</b>  |
| <b>3. Data Exploration.....</b>               | <b>4</b>  |
| <b>4. Data Preparation .....</b>              | <b>6</b>  |
| <b>5. Model Development .....</b>             | <b>7</b>  |
| 5.1 Decision Trees .....                      | 7         |
| 5.2 Regressions .....                         | 14        |
| 5.3 Neural Networks.....                      | 17        |
| 5.4 Model Assessment.....                     | 22        |
| <b>6. Recommendations &amp; Insights.....</b> | <b>23</b> |
| 6.1. Key Insights .....                       | 23        |
| 6.2. Model-Specific Insights.....             | 24        |
| 6.3. Actionable Recommendations .....         | 25        |
| 6.3.1. Offer Diverse Product Sets .....       | 25        |
| Recommendation.....                           | 25        |
| Strategies to Implement .....                 | 25        |
| 6.3.2. Gender-Specific Strategies .....       | 26        |
| Recommendation.....                           | 26        |
| Strategies to Implement .....                 | 26        |
| 6.3.3. Customer Engagement Initiatives.....   | 27        |
| Recommendation.....                           | 27        |
| Strategies to Implement .....                 | 27        |
| 6.3.4. Pilot Programs and A/B Testing .....   | 28        |
| Recommendation.....                           | 28        |
| Strategies to Implement .....                 | 28        |
| 6.3.5. Product Bundling.....                  | 29        |
| Recommendation.....                           | 29        |
| Strategies to Implement .....                 | 29        |
| 6.3.6. Data-Driven Monitoring .....           | 30        |
| Recommendation.....                           | 30        |
| Strategies to Implement .....                 | 30        |
| 6.4. Strategic Implementation Plan .....      | 31        |
| <b>7. Conclusion.....</b>                     | <b>31</b> |



## Dataset: RBC Churn Dataset

### 1. Executive Summary

Customer churn presents a significant challenge for RBC, impacting revenue and growth. Bailyn (2024) highlights that customer acquisition costs (CAC) in banking vary significantly, with retail banks averaging \$561 per customer, emphasizing the importance of cost-effective retention strategies to sustain profitability. These high acquisition costs underscore the importance of customer retention as a cost-effective strategy to sustain profitability. This report focuses on predicting customer churn for RBC using a dataset of customer demographics, account details, and behavioral attributes. The goal was to identify patterns and factors contributing to churn and to develop predictive models that accurately forecast whether a customer would stay with or leave the bank.

The project involved several key steps:

- **Data Exploration:**

The dataset was analyzed to understand the distribution and characteristics of key variables such as age, balance, credit score, and tenure. Initial data exploration highlighted skewness in two variables: age and number of products.

- **Data Preparation:**

Transformations (i.e., Cap and Floor and log transformations) were applied to reduce skewness in variables. The log transformation reduced skewness for Age. New binary variables like **HasBalance** were engineered to enhance model performance. The variable **HasProducts** was added to fix the skewness of the variable **Number of Products**. The dataset was partitioned into 50% training and 50% validation sets.

- **Model Development:**

Three types of predictive models were built, and the average squared error (ASE) were compared to determine the most accurate model:

- **Decision Tree:** A simple and interpretable model with an ASE of 0.124137.
- **Logistic Regression:** Provided insights into key predictors with an ASE of 0.129104.
- **Neural Network (Cap and Floor):** Achieved the best performance with an ASE of 0.01755 at 60 iterations and a ROC of 0.115123.

- **Model Comparison:**

**The Neural Network (Cap & Floor) model achieved the best performance with the highest ROC of 0.833 and second lowest ASE of 0.115865.**



- **Key Findings:**

- **Factors Influencing Churn:** Gender, fewer products, inactive membership, and certain geographic regions (e.g., Germany) were significant predictors of churn.
- **Churn Rate:** The overall churn rate in the dataset was **[20.37%]** as seen in the dataset.

- **Recommendations:**

- Implement targeted retention strategies for inactive members and those with low credit scores.
- Develop region-specific outreach programs to address higher churn rates in specific areas.
- Enhance customer engagement through personalized services and loyalty programs.

This analysis provides actionable insights for RBC to reduce customer churn and improve retention strategies. The neural network model serves as a robust tool for predicting at-risk customers, helping RBC take proactive measures to retain them.

## Introduction

Customer churn poses a significant challenge for the banking industry. Studies show that financial institutions lose approximately \$1 trillion annually due to customer attrition. Furthermore, acquiring new customers can cost up to 10 times more than retaining existing customers. This highlights the critical need for effective churn prediction and proactive retention strategies.

The goal of this project was to analyze RBC's customer data to identify patterns and predictors of churn. By developing robust predictive models, we aim to provide RBC with actionable insights to reduce churn rates and foster long-term customer relationships.

The dataset includes demographic details, account information, and behavioral metrics. The target variable, **Exited**, indicates whether a customer has left the bank (1) or remained (0). Our analysis focuses on understanding the drivers of churn and recommending strategies to retain high-risk customers.



## 2. Objective

The objective of this project is to develop predictive models to estimate which customers will churn for RBC. The goal is to identify key factors that influence whether a customer will stay or leave the bank. By analyzing customer demographics, account details, and behavioral attributes, the project aims to:

1. **Build and compare predictive models** (Decision Tree, Logistic Regression, and Neural Network) to forecast churn accurately.
2. **Identify significant predictors** of customer churn.
3. **Provide actionable insights** and recommendations to help RBC reduce churn rates and improve customer retention strategies.

### Dataset Source

- **Dataset:** RBC Churn Dataset
- **Source:** <https://www.kaggle.com/datasets/saadsalim997/rbcchurndataset>



### 3. Data Exploration

#### Variables Overview

After obtaining the dataset, an exploratory analysis was conducted to understand its structure. The dataset comprises **10,000 customer records** with the following key variables:

| Variable         | Description   | Type     |
|------------------|---|----------|
| Age              | Age of customer   | Interval |
| Balance          | Customer bank balance   | Interval |
| Credit Score     | Customer's credit score                                       | Interval |
| CustomerId       | Unique ID number of customer                                  | ID       |
| Estimated Salary | Approx. salary for each customer                              | Interval |
| GenderID         | Female (2), Male (1)  | Binary   |
| GeographyID      | France (1), Spain (2), Germany (3)                            | Nominal  |
| HasBalance       | Binary indicator for bank balance [Yes (1) / No (0)]          | Binary   |
| HasCrCard        | Binary Indicator for credit card ownership [Yes (1) / No (0)] | Binary   |
| IsActiveMember   | Active (1), Inactive (0)                                      | Binary   |
| Tenure           | Number of years that the customer has been with the bank      | Interval |
| Has Products     | 1 product (0), multiple products (1)                          | Nominal  |

#### Rejected Variables

- **Bank DOJ:** Rejected as it is highly correlated with *Tenure* which is an easier variable to model.
- **RowNumber:** Rejected due to lack of relevance for prediction.
- **NumofProducts:** Rejected and used *HasProducts* which has more significance.

#### Dataset Summary

- **Total Records:** 160,000 records
- **Missing Values:** 0 missing values
- **Variable Types:**
  - **Interval Variables:** *Age, Balance, Credit Score, Estimated Salary, Tenure.*
  - **Binary Variables:** *Exited, GenderID, HasBalance, HasCrCard, IsActiveMember.*
  - **Nominal Variables:** *CustomerId, GeographyID, HasProducts*



RBC Churn(2)

Variables - FIMPORT

(none) ☐ not Equal to ☐ ...

Columns: ☐ Label ☐ Mining ☐ Basic ☐ Statistics

| Name            | Role     | Level    | Report | Order | Drop | Lower Limit | Upper Limit |
|-----------------|----------|----------|--------|-------|------|-------------|-------------|
| Age             | Input    | Interval | No     |       | No   | .           | .           |
| Balance         | Input    | Interval | No     |       | No   | .           | .           |
| Bank_DOJ        | Rejected | Interval | No     |       | No   | .           | .           |
| CreditScore     | Input    | Interval | No     |       | No   | .           | .           |
| CustomerId      | ID       | Nominal  | No     |       | No   | .           | .           |
| EstimatedSalary | Input    | Interval | No     |       | No   | .           | .           |
| Exited          | Target   | Binary   | No     |       | No   | .           | .           |
| GenderID        | Input    | Binary   | No     |       | No   | .           | .           |
| GeographyID     | Input    | Nominal  | No     |       | No   | .           | .           |
| HasBalance      | Input    | Binary   | No     |       | No   | .           | .           |
| HasCrCard       | Input    | Binary   | No     |       | No   | .           | .           |
| HasProducts     | Input    | Nominal  | No     |       | No   | .           | .           |
| IsActiveMember  | Input    | Binary   | No     |       | No   | .           | .           |
| NumOfProducts   | Rejected | Interval | No     |       | No   | .           | .           |
| RowNumber       | Rejected | Interval | No     |       | No   | .           | .           |
| Tenure          | Input    | Interval | No     |       | No   | .           | .           |

Figure 1

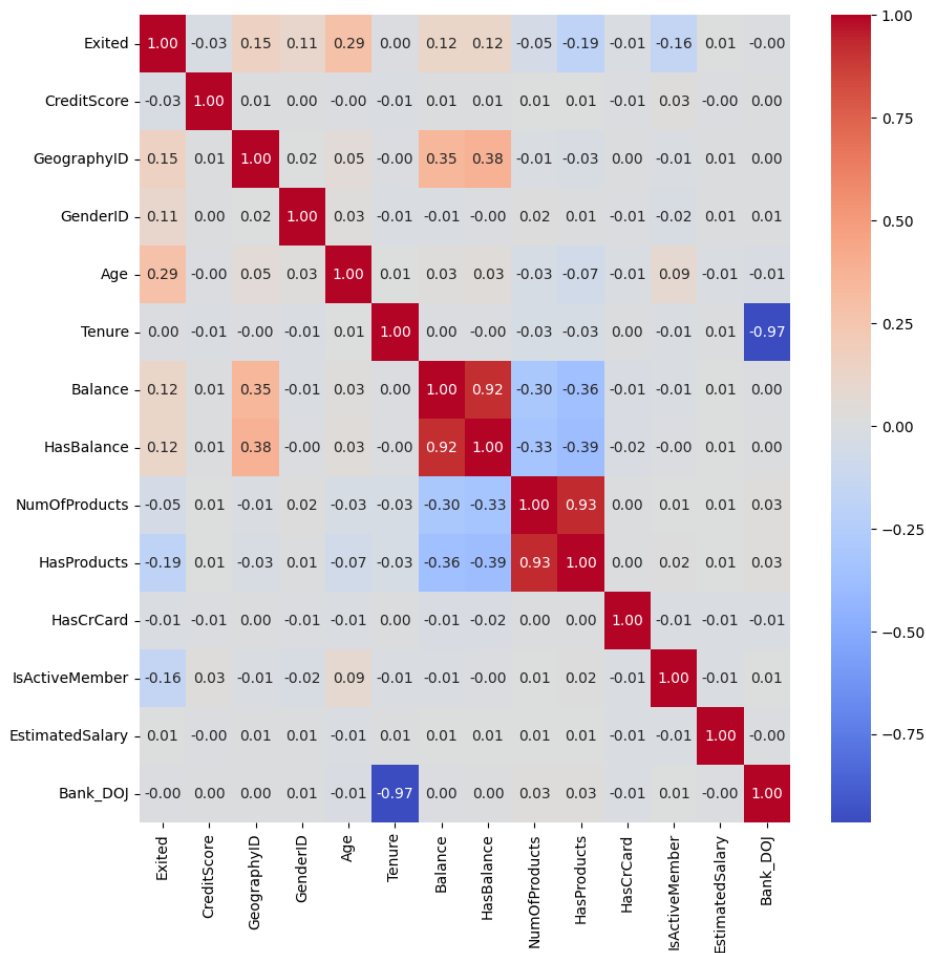


Figure 2



## 4. Data Preparation

### Steps Taken

1. Rejected Variables:
  - **RowNumber**, and **NumofProducts** were excluded as they were deemed irrelevant.
  - A heatmap was created to determine the correlation between our target variable 'Exited' and the other features. **Bank\_DOJ** had no correlation with the target variable; therefore, it was decided to exclude this feature. Additionally, **Bank\_DOJ** was a raw date field that could not be easily transformed into a useful variable, since calculating account tenure would duplicate information already captured by the Tenure variable.
2. Handling Missing Values and Data Validation
  - No missing values were included in the dataset
  - Variable ranges were checked to rule out unrealistic values (e.g., extreme ages and tenure, negative balances, invalid tenure)
3. Data Partition:
  - The data was split into **50% Train** and **50% Validation** sets.
4. Transformations:
  - Impute Variables: There were no missing or null values for variables in this dataset, therefore imputation was not required.
  - Skewed Variables:
    - **Cap and Floor** was applied to reduce skewness in the **Age** variable. Skewness was reduced but remained. (Figure 3)
    - **Log Transformation** was applied to further reduce skewness in the **Age** variable. After transformation, the **Age** variable was no longer skewed.
5. Feature Engineering (creating new feature):
  - A binary feature **HasProducts** was created as an indicator of customers with one product (0) and more than one product (1)
  - Another binary feature, **HasBalance**, was created as an indicator of customers who have no balance (0) or a positive balance in their account (1)
6. Dimensionality Reduction
  - Although dimensionality reduction techniques (e.g., collapsing categories) were considered, the number of variables was sufficiently small so no further dimensionality reduction was needed.





## 5. Model Development

### 5.1 Decision Trees

Three types of decision tree models were developed to predict customer churn: **Assessment Tree**, **Maximal Tree**, and **Average Squared Error (ASE) Tree**. Each tree was configured with specific settings and evaluated based on their **ASE** and other fit statistics.

#### Assessment Tree

The subtree method was set to “Assessment” and the assessment measure to “Decision”. After running the model, the ASE was 0.130702. There were 9 leaves, with **Age** as the first split and competing splits of **IsActiveMember**, **GeographyID** and **HasProduct**.

#### Configuration:

- **Subtree Method:** Assessment
- **Assessment Measure:** Decision
- **First Split:** Age
- **Competing Splits:** IsActiveMember”, “GeographyID” and “HasProduct”

#### Fit Statistics:

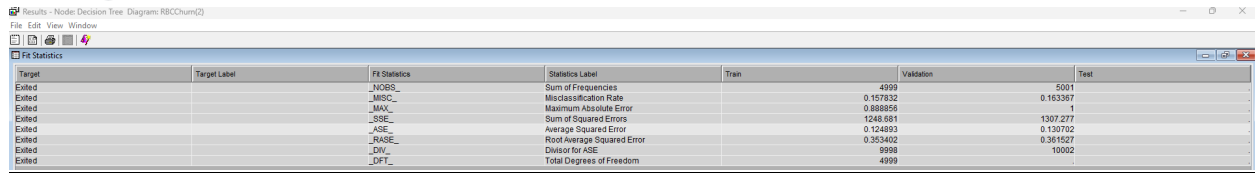
| Statistic                   | Train    | Validation |
|-----------------------------|----------|------------|
| Misclassification Rate      | 0.157832 | 0.163367   |
| Average Squared Error (ASE) | 0.124893 | 0.130702   |
| Root ASE                    | 0.353402 | 0.361527   |

#### Tree Details:

- Number of Leaves: 9
- First Split: Age
- Key Variables:
  - IsActiveMember
  - HasProduct
  - GeographyID

#### Insights:

This tree balances simplicity and performance, providing a clear understanding of how churn is influenced by age, product count, and balance status.



Results - Node Decision Tree Diagram: R8C0um2

File Edit View Window

Tree

```

graph TD
    Node1["Node Id: 1  
Statistic: 0  
Train: 75.00%  
Validation: 75.00%  
1: 20.36%  
2: 79.64%"]
    Node1 -->|Age| Node2["Node Id: 2  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node1 -->|Age| Node3["Node Id: 3  
Statistic: 0  
Train: 55.00%  
Validation: 60.00%  
1: 40.00%  
2: 60.00%"]
    Node2 -->|< 41.5 Or Missing| Node4["Node Id: 4  
Statistic: 0  
Train: 76.10%  
Validation: 76.10%  
1: 23.90%  
2: 76.10%"]
    Node2 -->|>= 41.5| Node5["Node Id: 5  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node3 -->|IsActiveMember| Node6["Node Id: 6  
Statistic: 0  
Train: 79.70%  
Validation: 79.70%  
1: 20.30%  
2: 79.70%"]
    Node3 -->|IsActiveMember| Node7["Node Id: 7  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node4 -->|GeographyID| Node8["Node Id: 8  
Statistic: 0  
Train: 65.00%  
Validation: 65.00%  
1: 35.00%  
2: 65.00%"]
    Node4 -->|GeographyID| Node9["Node Id: 9  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node5 -->|Balance| Node10["Node Id: 10  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node5 -->|Balance| Node11["Node Id: 11  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node6 -->|HasProducts| Node12["Node Id: 12  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node6 -->|HasProducts| Node13["Node Id: 13  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node7 -->|HasProducts| Node14["Node Id: 14  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node7 -->|HasProducts| Node15["Node Id: 15  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node8 -->|GeographyID| Node16["Node Id: 16  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node8 -->|GeographyID| Node17["Node Id: 17  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node9 -->|GeographyID| Node18["Node Id: 18  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node9 -->|GeographyID| Node19["Node Id: 19  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node10 -->|HasBalance| Node20["Node Id: 20  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node10 -->|HasBalance| Node21["Node Id: 21  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node11 -->|HasBalance| Node22["Node Id: 22  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node11 -->|HasBalance| Node23["Node Id: 23  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node12 -->|HasBalance| Node24["Node Id: 24  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node12 -->|HasBalance| Node25["Node Id: 25  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node13 -->|HasBalance| Node26["Node Id: 26  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node13 -->|HasBalance| Node27["Node Id: 27  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node14 -->|HasBalance| Node28["Node Id: 28  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node14 -->|HasBalance| Node29["Node Id: 29  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node15 -->|HasBalance| Node30["Node Id: 30  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node15 -->|HasBalance| Node31["Node Id: 31  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node16 -->|HasBalance| Node32["Node Id: 32  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node16 -->|HasBalance| Node33["Node Id: 33  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node17 -->|HasBalance| Node34["Node Id: 34  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node17 -->|HasBalance| Node35["Node Id: 35  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node18 -->|HasBalance| Node36["Node Id: 36  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node18 -->|HasBalance| Node37["Node Id: 37  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node19 -->|HasBalance| Node38["Node Id: 38  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node19 -->|HasBalance| Node39["Node Id: 39  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node20 -->|HasBalance| Node40["Node Id: 40  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node20 -->|HasBalance| Node41["Node Id: 41  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node21 -->|HasBalance| Node42["Node Id: 42  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node21 -->|HasBalance| Node43["Node Id: 43  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node22 -->|HasBalance| Node44["Node Id: 44  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node22 -->|HasBalance| Node45["Node Id: 45  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node23 -->|HasBalance| Node46["Node Id: 46  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node23 -->|HasBalance| Node47["Node Id: 47  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node24 -->|HasBalance| Node48["Node Id: 48  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node24 -->|HasBalance| Node49["Node Id: 49  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node25 -->|HasBalance| Node50["Node Id: 50  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node25 -->|HasBalance| Node51["Node Id: 51  
Statistic: 0  
Train: 88.89%  
Validation: 88.89%  
1: 11.11%  
2: 88.89%"]
    Node26 -->|HasBalance| Node52["Node Id: 52
```

The screenshot shows the RStudio interface with a 'Subtree Assessment Plot' window open. The plot displays the Average Squared Error (ASE) on the y-axis (ranging from 0.12 to 0.16) against the Number of Leaves on the x-axis (ranging from 0 to 30). Two lines are plotted: a blue line for 'Train Average Squared Error' and a red line for 'Valid Average Squared Error'. Both lines show a decreasing trend as the number of leaves increases. A vertical blue line is drawn at approximately 9 leaves, labeled with a '9' in a box. The legend at the bottom indicates that the blue line represents 'Train Average Squared Error' and the red line represents 'Valid Average Squared Error'.



### Maximal Tree

The maximal tree was built to determine all possible splits. The subtree method was changed to “Largest” and the assessment measure to “Decision”. With this model, the ASE was 0.124772.

There were 30 leaves with **Age** as the first split, and competing splits being **IsActiveMember**, **GeographyID** and **HasProducts**.

### Configuration:

- **Subtree Method:** Largest
- **Assessment Measure:** Decision
- **First Split:** Age
- **Competing Splits:** IsActiveMember, GeographyID, HasProducts, Balance, HasBalance and GenderID.

### Fit Statistics:

| Statistic                   | Train    | Validation |
|-----------------------------|----------|------------|
| Misclassification Rate      | 0.15103  | 0.163967   |
| Average Squared Error (ASE) | 0.115048 | 0.124772   |
| Root ASE                    | 0.339187 | 0.35323    |

### Tree Details:

- Number of Leaves: 30
- First Split: Age
- Key Variables:
  - IsActiveMember
  - GeographyID
  - HasProducts

### Insights:

The maximal tree captures all possible splits, providing detailed insights but potentially overfitting the data. It highlights complex interactions among variables.



Results - Node Maximal Tree Diagram RBCChurn(2)

File Edit View Window

Fit Statistics

| Target | Target Label | Fit Statistics | Statistics Label           | Train | Validation | Test     |
|--------|--------------|----------------|----------------------------|-------|------------|----------|
| Edited |              | _NOBS_         | Sum of Frequencies         |       | 4999       | 5001     |
| Edited |              | _MISC_         | Maximum Absolute Error     |       | 0.15103    | 0.163967 |
| Edited |              | _MAE_          | Maximum Absolute Error     |       | 0.071063   | 1        |
| Edited |              | _SSE_          | Sum of Squared Errors      |       | 1150.251   | 1247.966 |
| Edited |              | _ASE_          | Average Squared Error      |       | 0.115048   | 0.124772 |
| Edited |              | _RASE_         | Root Average Squared Error |       | 0.339167   | 0.35323  |
| Edited |              | _DIV_          | Divisor for ASE            |       | 9998       | 10002    |
| Edited |              | _DFT_          | Total Degrees of Freedom   |       | 4999       |          |

Figure 5

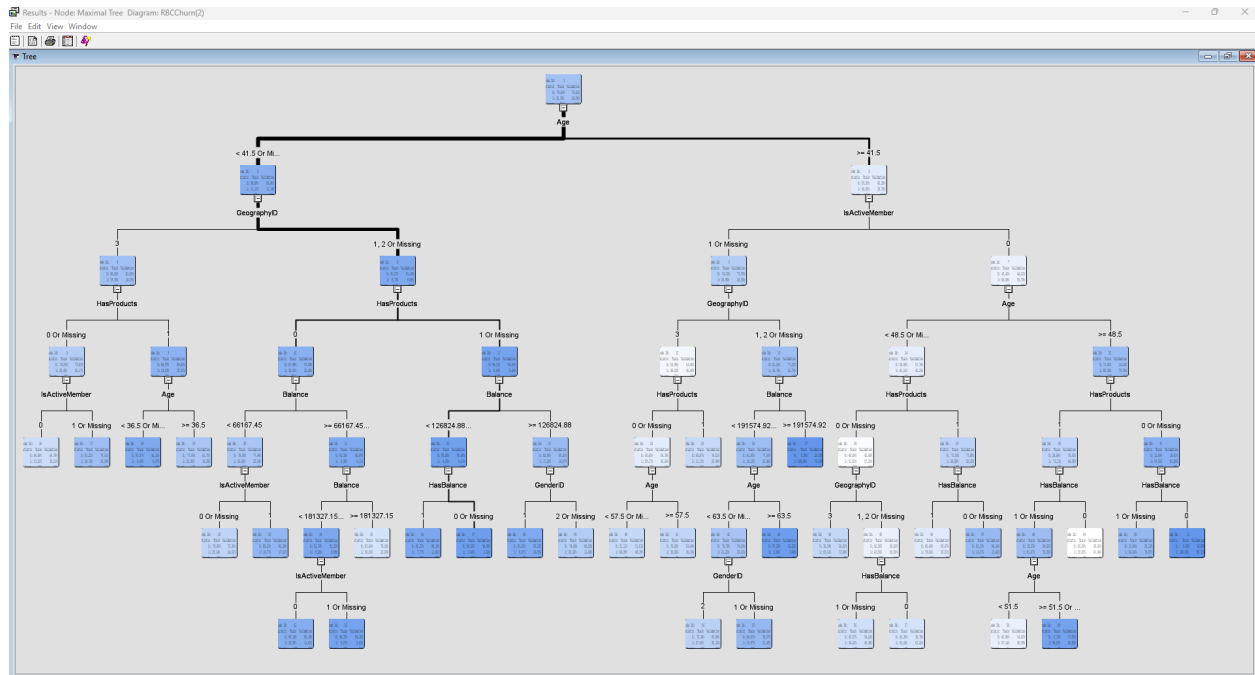


Figure 6

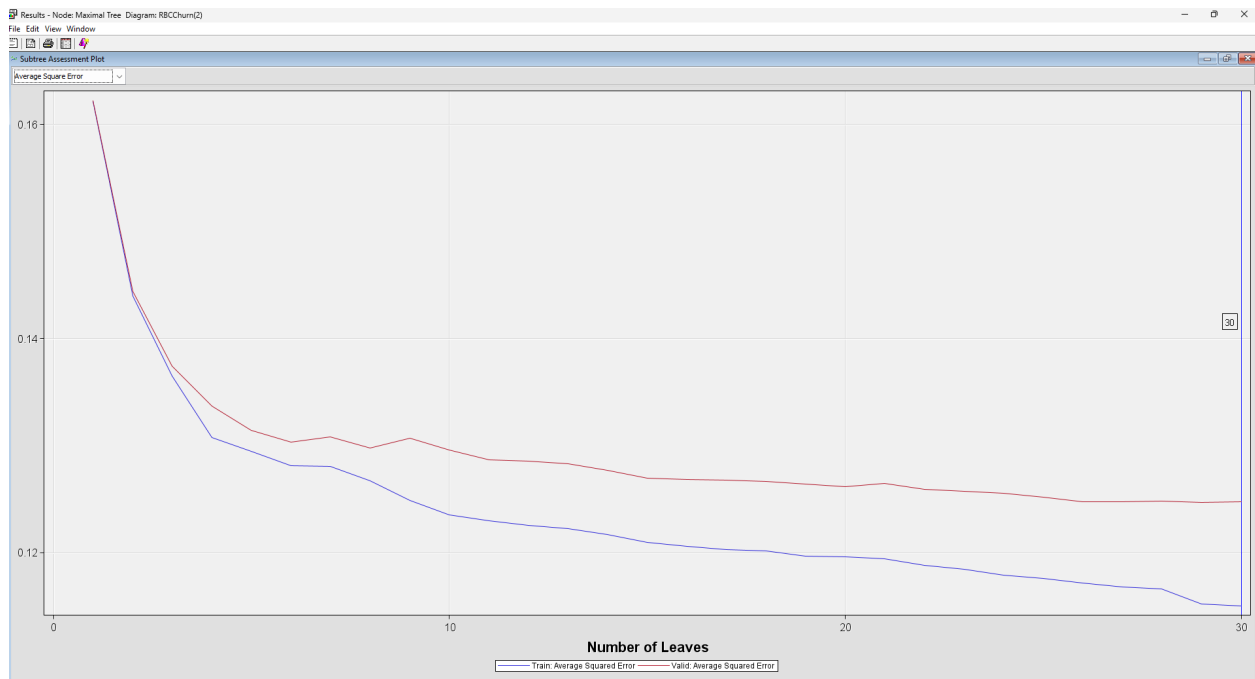


Figure 7



## ASE Tree

An ASE tree was built to optimize the best split. The subtree method was changed to “Assessment” and the assessment measure to “Average Square Error”. After running this model, the ASE was 0.124137. There were 26 leaves with “Age” as the first split and competing splits of **IsActiveMember**, **GeographyID**, **HasProducts**, **Balance**, **HasBalance** and **GenderID**.

### Configuration:

- **Subtree Method:** Assessment
- **Assessment Measure:** Average Square Error
- **First Split:** Age
- **Competing Splits:** IsActiveMember, GeographyID, HasProducts

### Fit Statistics:

| Statistic                   | Train    | Validation |
|-----------------------------|----------|------------|
| Misclassification Rate      | 0.153431 | 0.165767   |
| Average Squared Error (ASE) | 0.116497 | 0.124137   |
| Root ASE                    | 0.341316 | 0.352331   |

### Tree Details:

- Number of Leaves: 26
- First Split: Age
- Key Variables:
  - IsActiveMember
  - GeographyID
  - HasProducts

### Insights:

The ASE tree was optimized to minimize the Average Squared Error. It offers a balance between detail and performance and avoids excessive overfitting.



Results - Node ASE Tree Diagram RBCChurn(2)

File Edit View Window

Fit Statistics

| Target | Target Label | Fit Statistics | Statistics Label           | Train    | Validation | Test |
|--------|--------------|----------------|----------------------------|----------|------------|------|
| Exited |              | _NOBS_         | Sum of Frequencies         | 4999     | 5001       |      |
| Exited |              | _MSE_          | Maximum Absolute Error     | 0.153431 | 0.165767   |      |
| Exited |              | _MAE_          | Maximum Absolute Error     | 0.071953 | 1          |      |
| Exited |              | _SSE_          | Sum of Squared Errors      | 1164.735 | 1241.619   |      |
| Exited |              | _ASE_          | Average Squared Error      | 0.116457 | 0.124137   |      |
| Exited |              | _RASE_         | Root Average Squared Error | 0.341316 | 0.352351   |      |
| Exited |              | _DIV_          | Divisor for ASE            | 9998     | 10002      |      |
| Exited |              | _DFT_          | Total Degrees of Freedom   | 4999     |            |      |

Figure 8

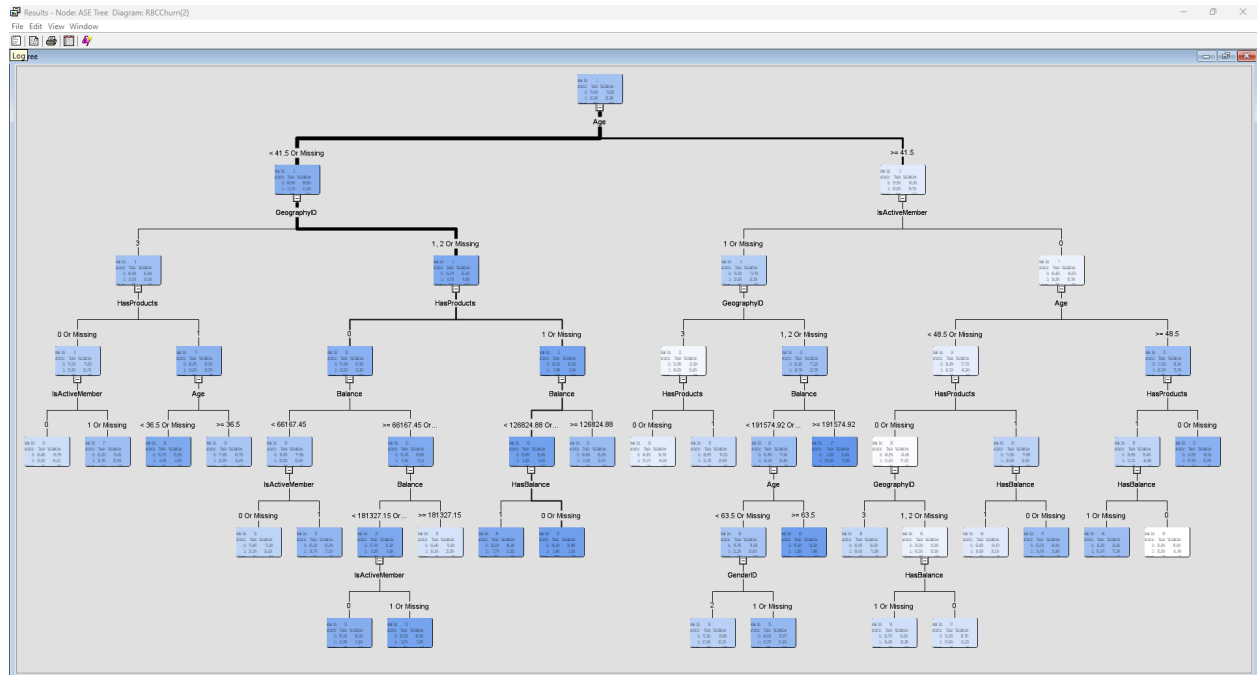


Figure 90

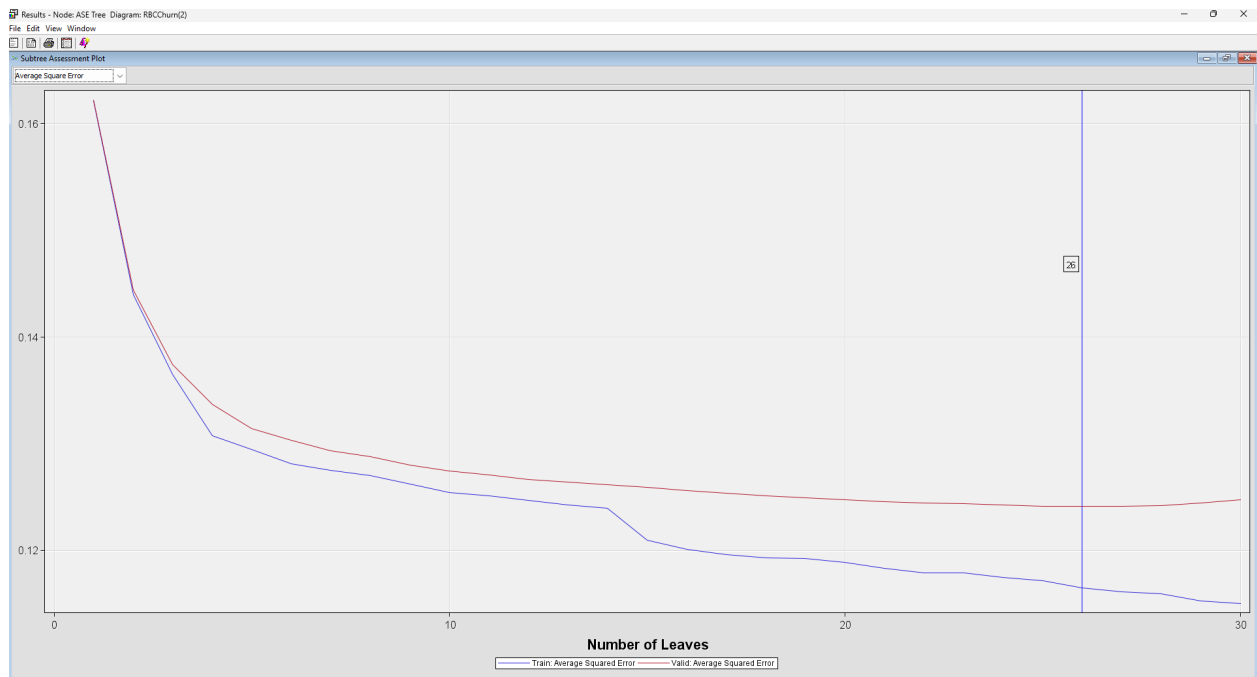


Figure 101



## Decision Tree Comparison

| Tree Type       | ASE      | Misclassification Rate | Number of Leaves | Key Splits                                    |
|-----------------|----------|------------------------|------------------|---|
| Assessment Tree | 0.130702 | 0.163367               | 9                | Age, IsActiveMember, GeographyID, HasProducts |
| Maximal Tree    | 0.124772 | 0.163967               | 30               | Age, IsActiveMember, GeographyID              |
| ASE Tree        | 0.124137 | 0.165767               | 26               | Age, IsActiveMember, GeographyID              |

## Key insights on Trees

The analysis across the ASE, Maximal, and Assessment Trees consistently identifies age and geography (specifically France and Spain) as key factors influencing customer churn and financial behavior. Customers younger than 41.5 years in France and Spain tend to have lower account balances and exhibit lower activity levels, making them more prone to churn.

Conversely, customers older than 41.5 years generally maintain higher balances and display greater engagement, such as owning more financial products and demonstrating higher activity levels. Notably, within the older segment, those over 48.5 years with only one product tend to show reduced financial activity, indicating another distinct subgroup with unique retention challenges.

## Conclusion

- The ASE Tree provided the lowest ASE (0.124137), making it the most accurate decision tree model.
- The Maximal Tree offered the most detailed splits but risked overfitting.
- The Assessment Tree was simpler and interpretable but had a higher ASE.

These decision tree models highlight key factors influencing customer churn, such as **Age**, **IsActiveMember**, and **GeographyID**.



## 5.2 Regressions

After building a decision tree, regression models were developed (i.e., full regression, backward elimination, forward regression, and stepwise regression). These models helped refine the analysis, validate variable significance, and optimize the model for improved accuracy and interpretability

### Full Regression

For the full regression model, the ASE was 0.129104. The model had significant variables with p-values lower than 0.000. These variables were **GenderID**, **GeographyID**, **HasProducts**, **IsActiveMember**, and **Log\_REP\_Age**.

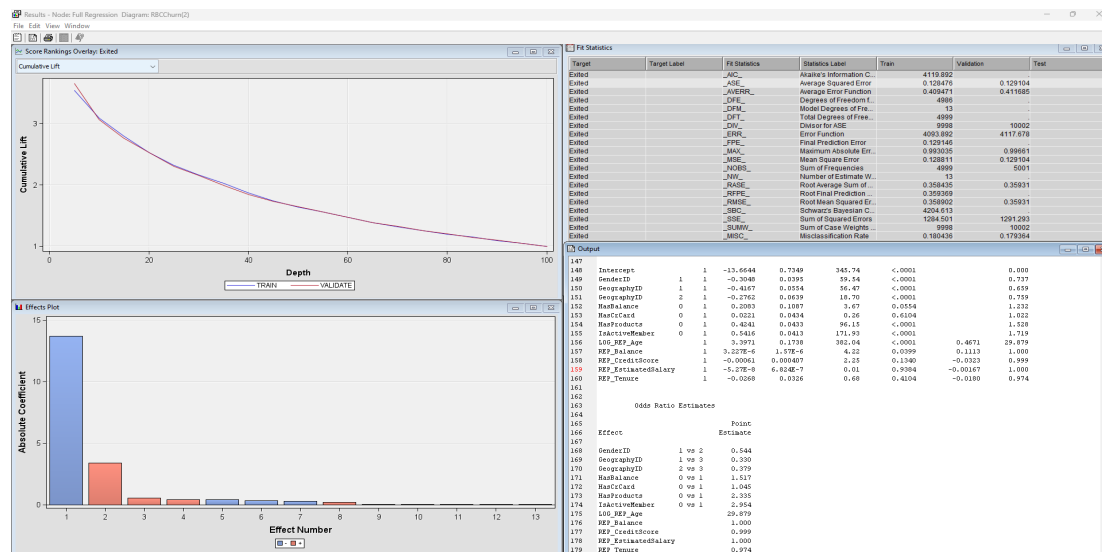


Figure 11

### Backward Regression

Backward regression was modeled to remove the least significant variable to optimize model complexity. The selection model was changed to “Backward” and the Selection Criterion to “Validation Error”. The ASE of the backward regression model was 0.12913.

Significant variables with p-values lower than 0.0001 were **GenderID**, **GeographyID**, **HasProducts**, **IsActiveMember** and **Log\_REP\_Age**.

Viewing the points estimate, it was determined how these significant variables influenced customer churn:

- GenderID” (1 vs 2) at 0.547, which indicates males are 45.3% less likely to churn compared to females.
- “GeographyID” (1 vs 3) at 0.334, which indicates that customers in France are 66.6% less likely to churn compared to customers in Germany.
- “GeographyID” (2 vs 3) at 0.384, which indicates that customers in Spain are 61.6% less likely to churn when compared to customers in Germany.
- “HasProducts” (0 vs 1) at 2.316, which indicates that people with one product are 131.6% more likely to churn than people with more than one product.





- “IsActiveMember” (0 vs 1) at 2.953, which indicates that customers who are inactive are nearly three times more likely to churn compared to customers that are active,
- “Log\_REP\_Age” (29.669) indicates that for every 1-unit increase in natural log of age (equivalent to a 2.718 -fold increase in age), customers are approximately 29.7 times more likely to churn. This suggests that older customers are more likely to churn than the younger ones.

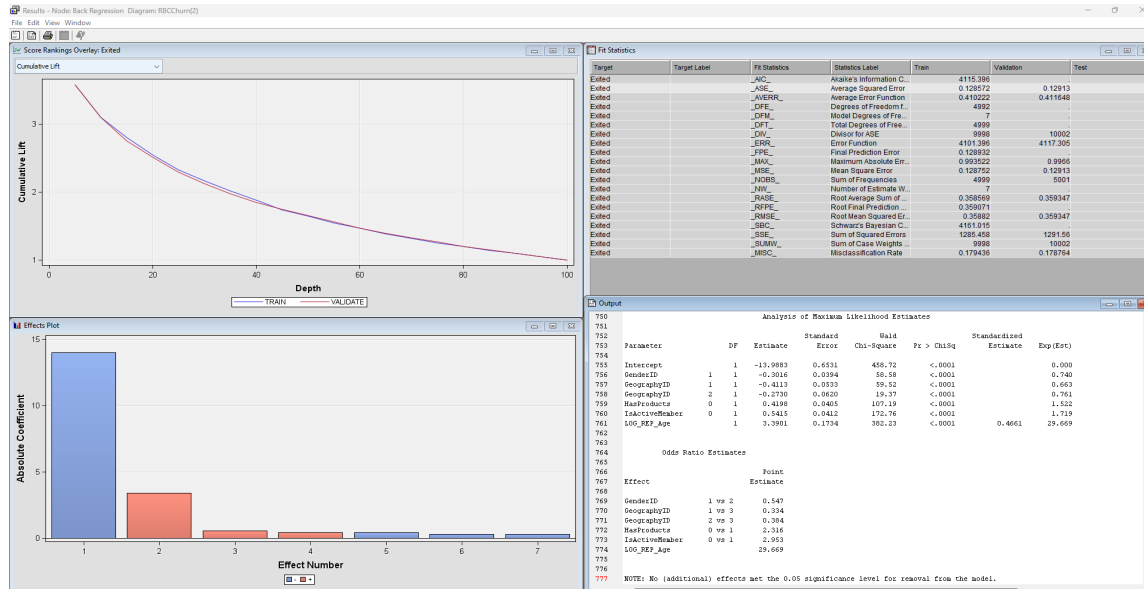


Figure 12

## Forward Regression

Forward regression was utilized to refine our model by sequentially adding variables based on their statistical significance to optimize model complexity. The selection model was changed to “Forward” and the Selection Criterion to “Validation Error” and ran the model. The ASE for forward regression was 0.12913.

Points estimates and p-values were the same as the backward regression model. For further analysis, please refer to Backward regression.

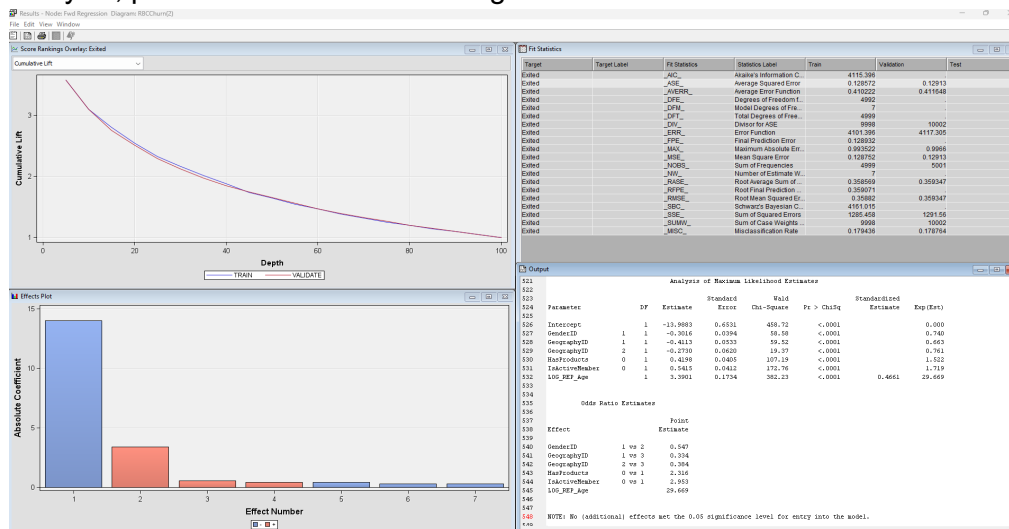




Figure 13

## Stepwise

Stepwise regression was used to iteratively add and remove variables based on their statistical significance and contribution to the model. This approach helped to identify the optimal set of predictors, balancing model complexity, and predictive accuracy.

For the stepwise model, the selection model was changed to “Stepwise” and the Selection Criterion to “Validation Error” and ran the model. The ASE of the stepwise model was 0.12913.

Points estimates and p-values were the same as the backward regression model. For further analysis, please refer to Backward regression.

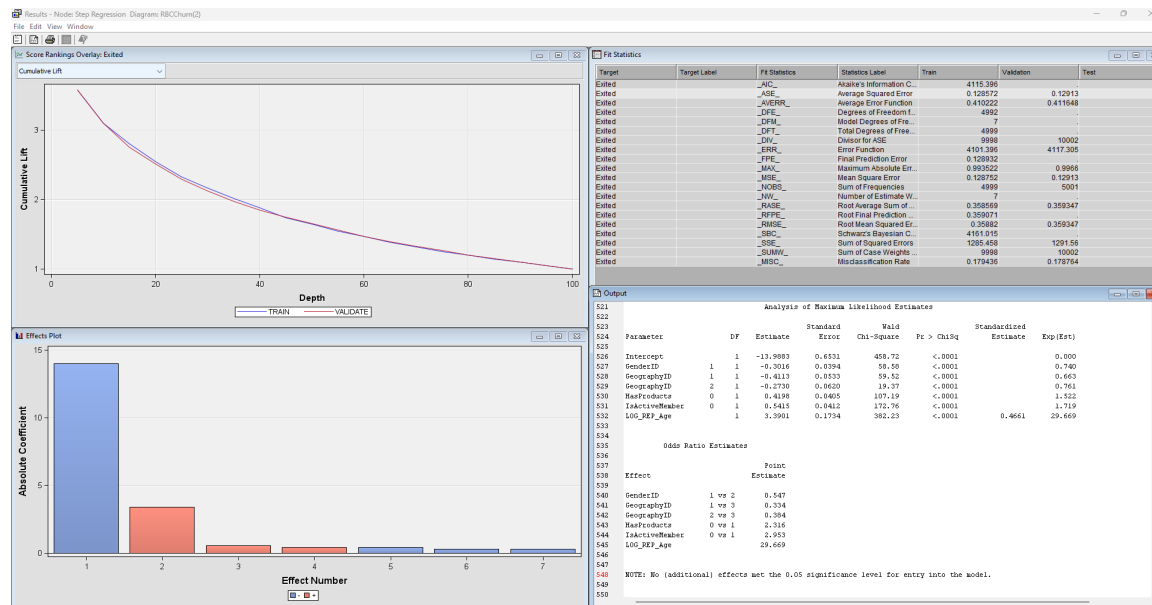


Figure 145



## 5.3 Neural Networks

Neural network models were built to capture complex non-linear relationships between variables and enhance the predictive accuracy of the analysis.

### Neural Network: Cap and Floor

One of the neural network models was connected to the cap and floor nodes. Model selection criterion was changed to “Average Error”, with Maximum iterations of “100”. Preliminary training was also disabled.

The model converged at 63 iterations, and the ASE was 0.115865.

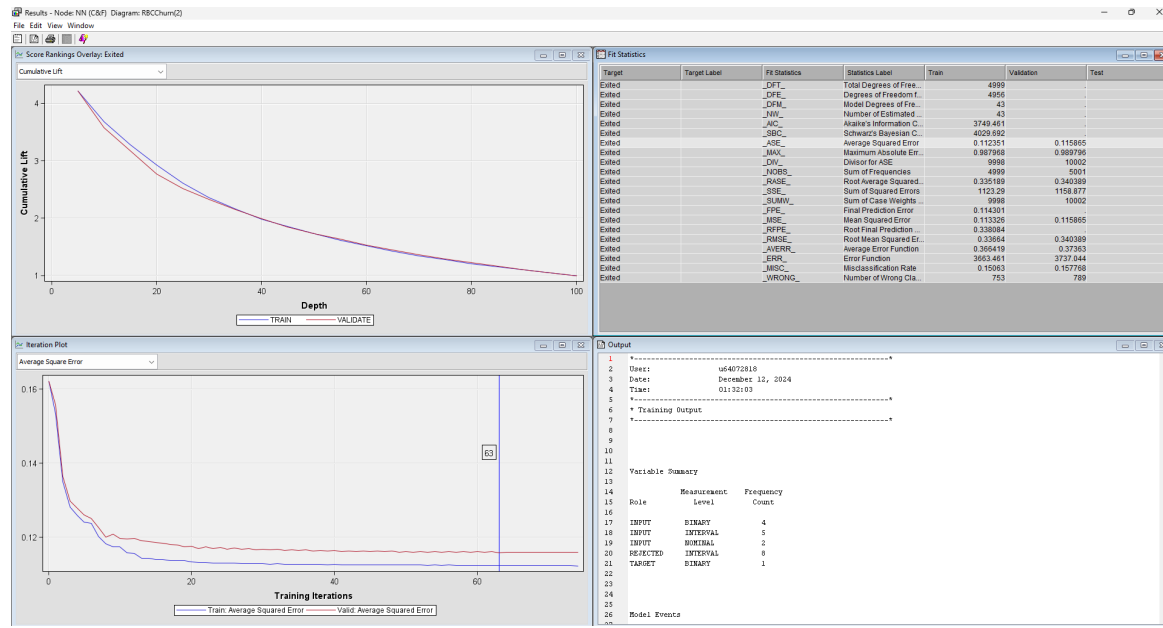


Figure 156



## Neural Network: Transform

The second neural network model was connected to the transform node. This model had 30 iterations, with an ASE of 0.121844.

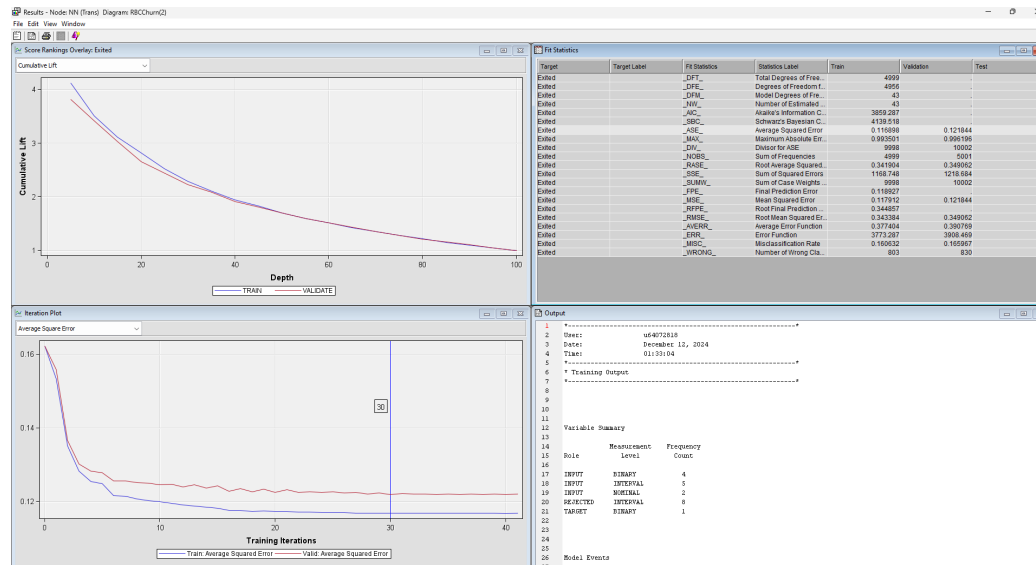


Figure 167

Since the Neural Network: Cap & Floor had a lower ASE, additional neural networks based on this model were run with different hidden units. By default, the neutral network runs with 3 hidden units.

## NN Cap and Floor 2H

The number of hidden units was changed to 2 to see if the model would perform better. The model had 48 iterations with an ASE of 0.119794. This was higher than the model with 3 hidden units.

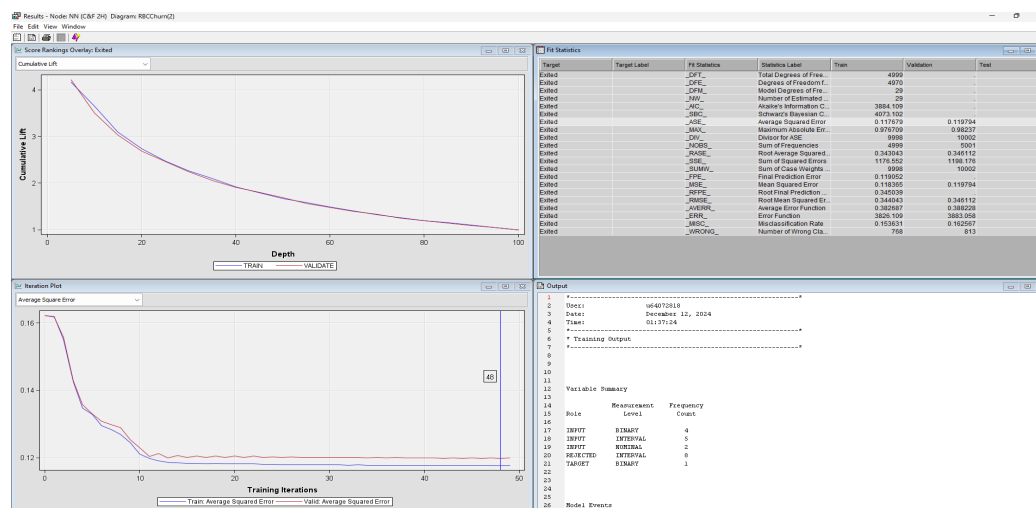


Figure 178



## NN Cap and Floor 4H

This model, with 4 hidden units, had 66 iterations with an ASE of 0.115123. This was lower than the models with 2 and 3 hidden units.

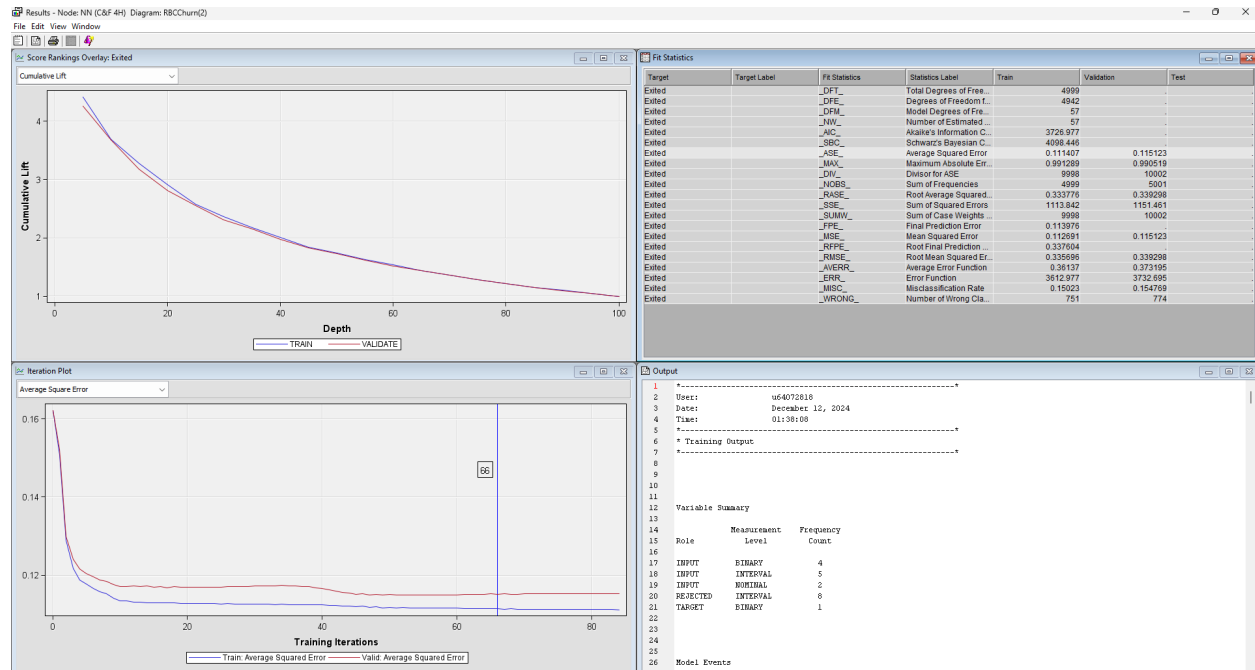


Figure 189

## NN Cap and Floor 5H

This model, with 5 hidden units, had 20 iterations with an ASE of 0.116445. This was higher than the Cap & Floor with 4 hidden units.

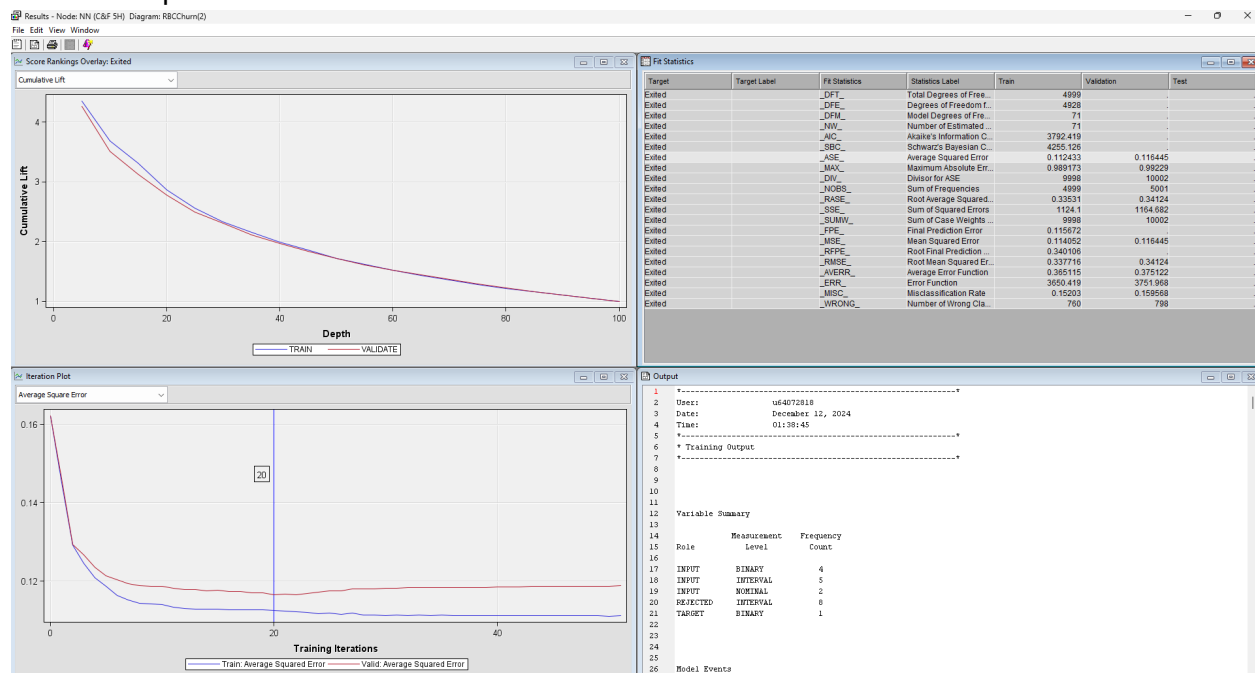




Figure 190

## NN Cap and Floor 6H

This model with 6 hidden units, had 11 iterations with an ASE of 0.116965. This was higher than the models with 3, 4, and 5 hidden units but lower than the model with 2 hidden units.

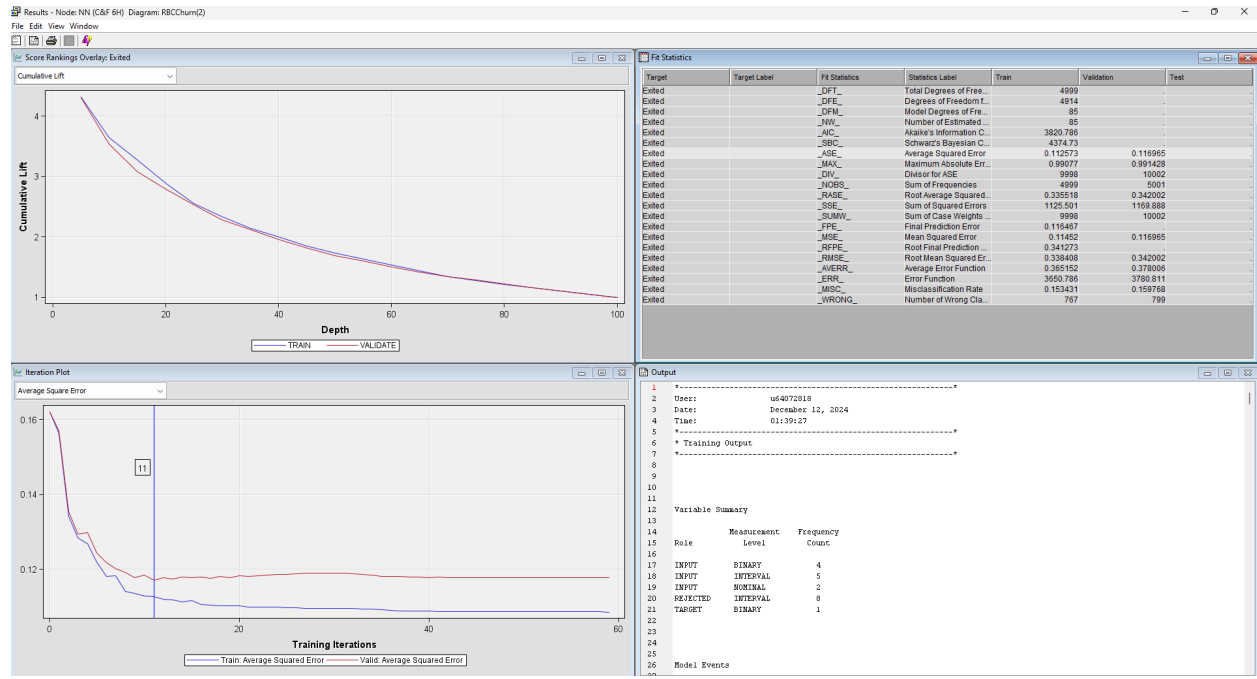


Figure 21

## NN Cap and Floor 7H

This model, with 7 hidden units, had 29 iterations with an ASE of 0.118729. This was the second highest ASE so far after the model with 2 hidden units.

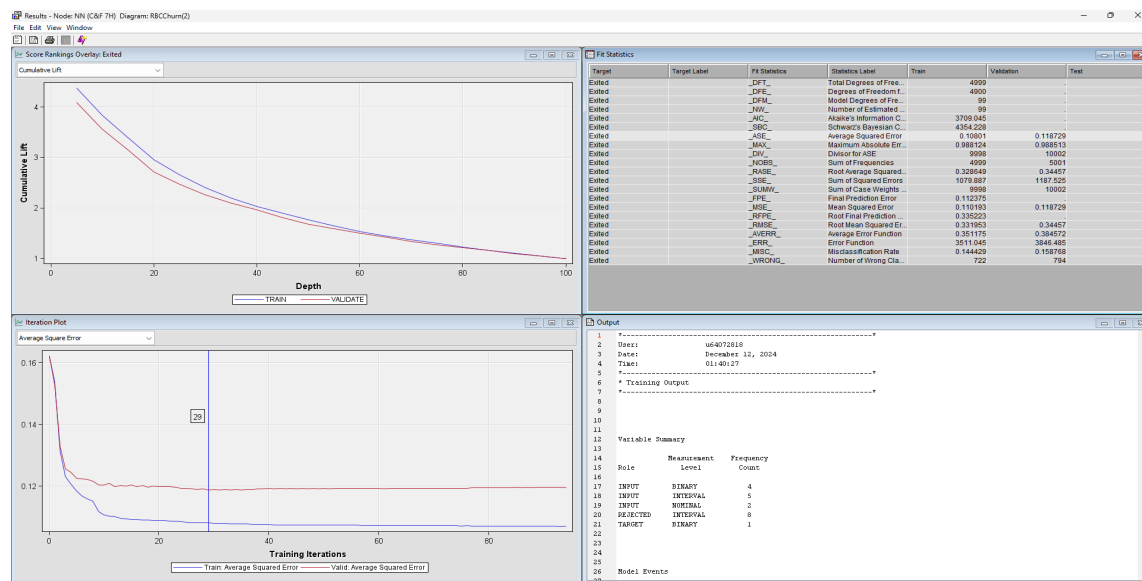


Figure 20



## NN Cap and Floor 8H

This model, with 8 hidden units, had 21 iterations with an ASE of 0.118568. This was the third highest ASE of the hidden unit Neural Network models.

No additional Neural network models were run after 8 hidden units as the ASE was increasing, therefore the accuracy of the models would likely worsen with each additional hidden unit.

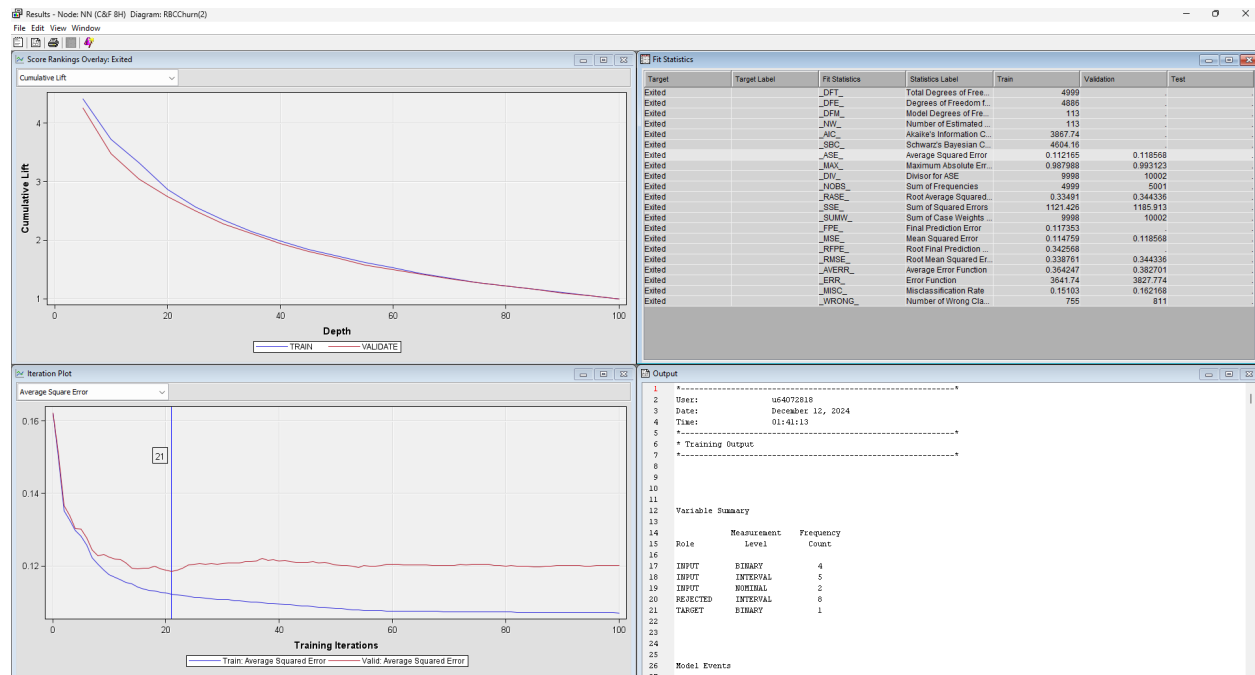


Figure 21

## Neural Network Summary Table

| Hidden units | ASE      | Number of Iterations |
|--------------|----------|----------------------|
| Cap & Floor  | 0.115865 | 63                   |
| Transform    | 0.12184  | 30                   |
| 2H           | 0.119794 | 48                   |
| 4H           | 0.115123 | 66                   |
| 5H           | 0.116445 | 20                   |
| 6H           | 0.116965 | 11                   |
| 7H           | 0.118729 | 29                   |
| 8H           | 0.118568 | 21                   |



## 5.4 Model Assessment

For the model comparison, the selection statistic was based on the receiver operating characteristic curve (ROC) and the validation dataset split.

After running the Model Assessment node, the results indicated that the best model was Neural Network: Cap & Floor (3H) which had a ROC value of 0.833.

In the cumulative lift chart, the NN Cap & Floor (6H) model has the best response rate at a depth of 5, with results being 4.301616.

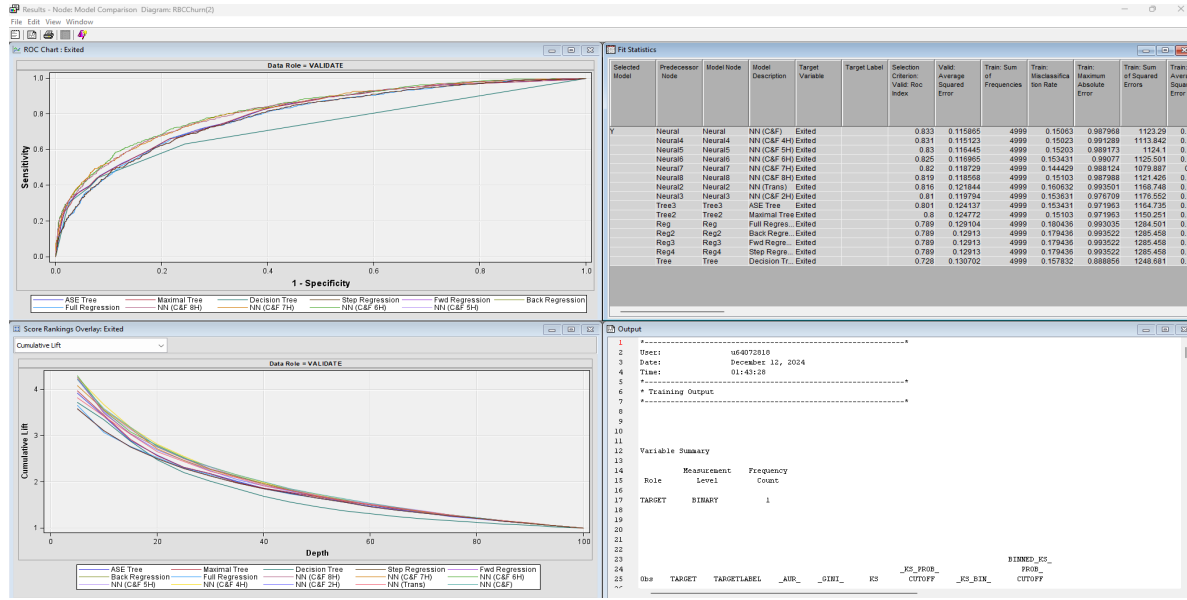


Figure 224

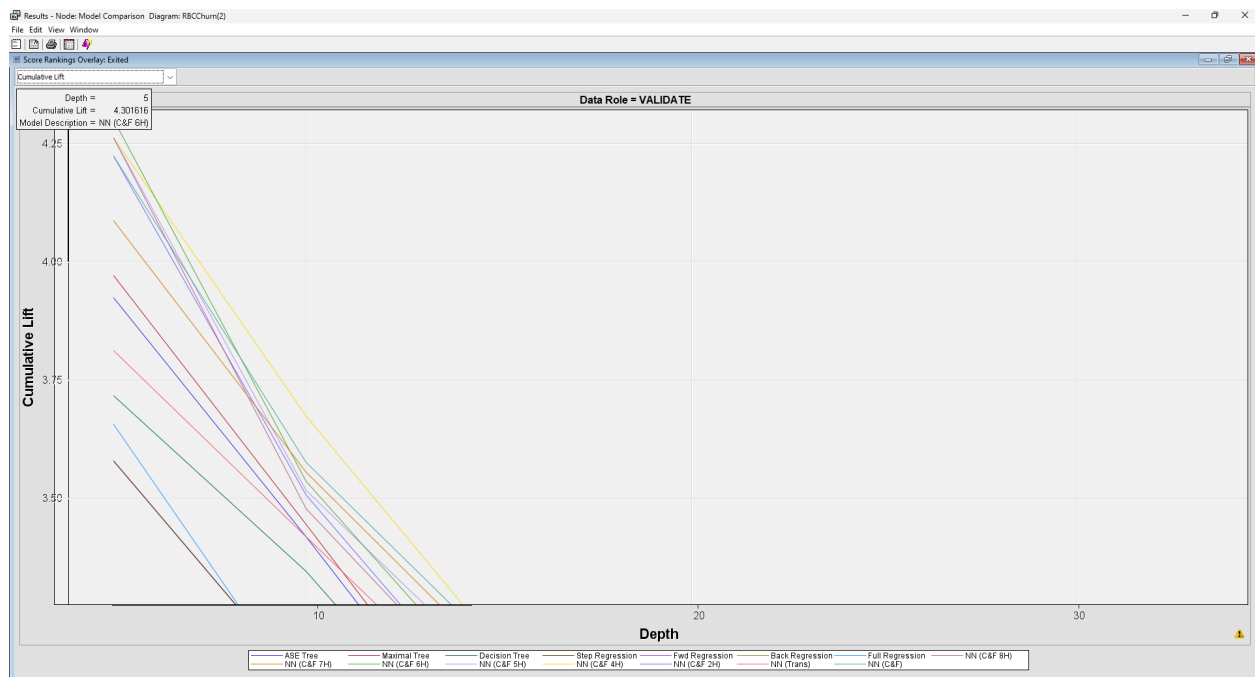


Figure 235





## 6. Recommendations & Insights

Based on the results of our models, several key insights have been identified as well as actionable recommendations to reduce customer churn and improve customer retention for RBC.

### 6.1. Key Insights

#### Demographic Factors

##### Gender:

- **Insight:** Males are **45.3% less likely** to churn compared to females.
- **Recommendation:** Develop targeted retention campaigns specifically aimed at female customers, addressing their needs and preferences.

##### Age:

- **Insight:** Older customers are more likely to churn compared to younger customers.
- **Recommendation:** Implement loyalty programs and personalized offers for older customers to increase engagement and satisfaction.

#### Geographic Factors

##### Geography:

- **Insight:** Customers in **France** and **Spain** are significantly less likely to churn compared to those in **Germany**.
- **Recommendations:**
  - Investigate potential reasons for higher churn rates in Germany (e.g., product offerings, customer service quality).
  - Develop region-specific strategies and incentives for German customers to improve retention.

#### Customer Behavior

##### Number of Products:

- **Insight:** Customers with only **one product** are significantly more likely to churn.
- **Recommendations:**
  - Encourage customers to adopt multiple products through bundling offers or cross-selling strategies.
  - Provide incentives or discounts for customers who add additional products.

#### Active Membership Status

- **Insight:** Inactive customers are about **2 times more likely** to churn compared to active members.
- **Recommendations:**
  - Launch engagement campaigns to reactivate dormant accounts.
  - Offer personalized incentives to encourage activity (e.g., cashback, special promotions).



## 6.2. Model-Specific Insights

### From Decision Trees

- Age is consistently the first and most important split in the decision trees.
- **Active Membership Status, GeographyID and HasProducts** are key factors influencing churn.

Recommendation: Focus retention efforts on younger customers with fewer products and inactive accounts.

### From Regression Models

- **Significant Predictors:**
  - **GenderID, GeographyID, HasProducts, IsActiveMember, and Log\_REP\_Age.**

### Recommendation:

- Target customers based on these predictors for personalized retention strategies.
- Use regression insights to understand customer segments and tailor communication accordingly.

### From Neural Networks

- **The Neural Network (Cap & Floor)** model achieved the best performance with the highest **ROC of 0.833** and second lowest **ASE of 0.115865**.
- The model identified complex non-linear relationships between features.

### Recommendations:

- Leverage the neural network insights to build advanced churn prediction tools.
- Use these tools for real-time churn risk assessments and proactive customer interventions.



## 6.3. Actionable Recommendations

### 6.3.1. Offer Diverse Product Sets

Older customers are at a higher risk of churning, often due to a lack of products that cater to their evolving financial needs. As customers age, their priorities shift towards retirement planning, investment growth, and wealth preservation, which may not be addressed by standard product offerings.

#### Recommendation

RBC should introduce personalized product bundles tailored specifically to older customers. These bundles should address their financial goals at different life stages, such as retirement, estate planning, and wealth management. By offering comprehensive solutions that evolve with their needs, RBC can foster long-term relationships and reduce churn.

#### Strategies to Implement

##### 1. Segment Analysis:

- Conduct a thorough analysis to identify different segments within the older customer demographic (e.g., pre-retirement, early retirement, and late retirement).
- Use data analytics to understand their financial behaviors, product preferences, and pain points.

##### 2. Personalized Bundling:

- Develop product bundles tailored to each segment, such as:
  - **“Retirement Ready Bundle”**: Includes retirement savings accounts, annuities, and estate planning services.
  - **“Investment Growth Plan”**: Combines investment portfolios, wealth advisory services, and tax-efficient investment strategies.
  - **“Legacy Protection Package”**: Focuses on wills, trusts, and life insurance products.

##### 3. Educational Campaigns:

- Launch educational initiatives like webinars, workshops, and personalized consultations to help older customers understand the benefits of these products.
- Develop resources (blogs, e-books, and videos) on topics such as retirement planning, wealth preservation, and legacy planning.

##### 4. Exclusive Incentives:

- Offer preferential rates on mortgages, reduced fees on investment services, and cashback for adopting bundled products.
- Create loyalty programs that reward older customers for engaging with multiple RBC products.

##### 5. Feedback Mechanisms:

- Regularly collect feedback through surveys and focus groups to refine and enhance these bundles based on customer needs.



### 6.3.2. Gender-Specific Strategies

Female customers are **42.5% more likely** to churn compared to male customers. This indicates a need for more targeted engagement strategies that address their unique financial needs, preferences, and experiences.

#### Recommendation

RBC should design and implement marketing campaigns and financial services tailored to female customers. These initiatives should focus on building trust, offering personalized solutions, and empowering women to achieve their financial goals.

#### Strategies to Implement

##### 1. Financial Empowerment Workshops:

- Host workshops and webinars specifically for women, covering topics such as:
  - Investment planning and portfolio management.
  - Budgeting, savings strategies, and debt management.
  - Retirement planning and financial independence.

##### 2. Personalized Communication:

- Use data-driven insights to create personalized communication strategies.
- Send tailored messages highlighting products and services that resonate with female customers, such as investment opportunities, savings plans, and family-oriented financial products.

##### 3. Success Stories and Case Studies:

- Share testimonials and case studies featuring successful female clients.
- Highlight how RBC's services have helped women achieve their financial goals, reinforcing RBC's commitment to supporting female customers.

##### 4. Dedicated Relationship Managers:

- Assign relationship managers who specialize in understanding and addressing the needs of female clients.
- Train staff to be sensitive to the unique financial challenges women may face, such as career breaks, wage gaps, and caregiving responsibilities.

##### 5. Women-Focused Financial Products:

- Develop products designed with women in mind, such as flexible investment plans, joint savings accounts for families, and financial planning services for working mothers.



### 6.3.3. Customer Engagement Initiatives

Inactive customers are twice as likely to churn compared to active customers. This highlights the importance of proactively re-engaging dormant customers to rekindle their interest in RBC's services.

#### **Recommendation**

RBC should develop targeted reactivation campaigns to encourage inactive customers to become active users again. Personalized incentives, timely communication, and value-driven offers can reignite customer engagement.

#### **Strategies to Implement**

##### **1. Cashback and Rewards Programs:**

- Introduce cashback offers or loyalty rewards for completing specific transactions, such as making deposits, applying for a credit card, or using online banking services.
- Create limited-time reward campaigns to create a sense of urgency and excitement.

##### **2. Exclusive Time-Sensitive Offers:**

- Provide promotions like discounted loan rates, waived fees, or bonus interest on savings accounts for a limited period.
- Tailor offers to each customer's historical behavior to maximize relevance.

##### **3. Personalized Outreach:**

- Use personalized emails, SMS, or phone calls to reach out to inactive customers.
- Craft messages that acknowledge their inactivity and present solutions or incentives to re-engage, such as "We've missed you! Here's an exclusive offer to welcome you back."

##### **4. Feedback Collection:**

- Reach out to inactive customers to understand the reasons behind their disengagement.
- Use surveys and one-on-one calls to gather insights and identify areas for service improvement.

##### **5. Re-Engagement Campaigns:**

- Launch campaigns with themes like "Welcome Back" or "Reconnect with RBC" to make customers feel valued and appreciated.
- Offer personalized product recommendations based on their previous engagement history.



#### 6.3.4. Pilot Programs and A/B Testing

Implementing new strategies at scale without testing can lead to inefficient use of resources and suboptimal outcomes. Testing ensures that only the most effective initiatives are rolled out widely.

##### **Recommendation**

RBC should conduct pilot programs and A/B testing to evaluate the effectiveness of new strategies before full-scale implementation. This method minimizes risk, optimizes resource allocation, and improves the likelihood of success.

##### **Strategies to Implement**

###### **1. Controlled Experiments:**

- Select a representative sample of customers to test new initiatives such as product bundles, marketing campaigns, or engagement tactics.
- Create control and test groups to compare outcomes objectively.

###### **2. Key Metrics Tracking:**

- Define and track critical metrics such as churn rate, customer engagement, product adoption rates, and return on investment (ROI).
- Analyze performance data to determine which strategies yield the best results.

###### **3. Iterative Refinements:**

- Based on the results of A/B tests, refine and adjust strategies to improve effectiveness.
- Implement a continuous improvement loop to ensure ongoing optimization.

###### **4. Documentation and Reporting:**

- Maintain detailed documentation of the testing process, results, and insights gained.
- Share reports with stakeholders to facilitate informed decision-making and ensure transparency.

###### **5. Scaling Successful Strategies:**

- Once a strategy proves successful in pilot testing, roll it out to the broader customer base with confidence.
- Develop implementation guidelines to ensure consistency and effectiveness across all regions.



### 6.3.5. Product Bundling

Customers who hold only a single product are significantly more likely to churn. Encouraging customers to adopt multiple products enhances their engagement, satisfaction, and loyalty, making them less likely to leave.

#### Recommendation

RBC should implement comprehensive product bundling strategies to incentivize customers to adopt multiple financial products. Bundling services such as **savings accounts, credit cards, mortgages, investment services, and insurance products** creates more value for customers and strengthens their relationship with RBC.

#### Strategies to Implement

##### 1. Service Bundling:

- Offer bundled packages that combine complementary services, such as:
  - **“Everyday Banking Bundle”**: Savings account + credit card + online banking.
  - **“Family Financial Bundle”**: Joint savings account + children’s education savings plan + family insurance.
  - **“Home Ownership Bundle”**: Mortgage + home insurance + home equity line of credit.

##### 2. Cross-Selling Campaigns:

- Launch targeted cross-selling campaigns that recommend additional products based on customers’ current holdings.
- Use personalized communication to highlight how adding products can benefit customers (e.g., “You already have a savings account! Add a credit card and enjoy cashback rewards on purchases.”).

##### 3. Incentives and Discounts:

- Offer discounts or perks for customers who adopt multiple products, such as:
  - **Fee Waivers**: Waive monthly fees for customers who bundle at least three services.
  - **Cashback Rewards**: Provide cashback on purchases made with bundled products.
  - **Interest Rate Benefits**: Offer lower mortgage rates or higher savings rates for bundled customers.

##### 4. Promotional Campaigns:

- Develop seasonal or limited-time promotions to encourage customers to adopt bundles (e.g., “Bundle now and receive a \$100 bonus!”).
- Highlight these promotions through emails, SMS, RBC’s website, and social media.

##### 5. Customer Education:

- Educate customers on the benefits of product bundling through workshops, webinars, and one-on-one consultations.
- Share success stories and testimonials from customers who have benefited from bundled services.

##### 6. Tracking and Personalization:

- Use customer data to identify which bundles are most relevant to specific segments.
- Implement AI-driven personalization to recommend the right bundle at the right time.



### 6.3.6. Data-Driven Monitoring

Predictive models, especially neural networks, provide powerful tools for identifying at-risk customers. Continuous monitoring of churn risk allows RBC to intervene in real time and proactively prevent customer attrition.

#### Recommendation

Leverage the capabilities of the **Neural Network model** to establish a robust data-driven monitoring system. This system should provide ongoing insights into churn risk, offer early warnings, and facilitate timely interventions to retain at-risk customers.

#### Strategies to Implement

##### 1. Continuous Risk Monitoring:

- Integrate the neural network model into RBC's customer relationship management (CRM) system to monitor churn risk continuously.
- Automate daily or weekly risk assessments to keep track of changes in customer behavior and churn likelihood.

##### 2. Early Warning System:

- Develop an early warning system that flags customers who exhibit high churn risk based on predictive scores.
- Set thresholds for churn probability (e.g., customers with a churn risk above **70%**) to trigger alerts for immediate follow-up by relationship managers.

##### 3. Real-Time Intervention Strategies:

- Implement real-time intervention strategies such as personalized offers, engagement calls, and targeted emails to retain high-risk customers.
- For example, if a customer shows signs of disengagement (e.g., reduced account activity), automatically send a personalized offer like "Enjoy a \$50 reward for your next transaction."

##### 4. Dynamic Customer Segmentation:

- Continuously update customer segments based on their churn risk scores.
- Create dynamic profiles that help RBC understand which segments are most vulnerable and tailor interventions accordingly.

##### 5. Dashboard and Reporting Tools:

- Develop interactive dashboards to visualize churn risk data, trends, and intervention outcomes.
- Provide regular reports to stakeholders highlighting key insights, intervention success rates, and areas for improvement.

##### 6. Feedback Loop:

- Establish a feedback mechanism to evaluate the effectiveness of interventions.
- Use insights from successful and unsuccessful interventions to continuously refine the monitoring and response processes.

##### 7. Integration with Customer Support:

- Equip customer support teams with churn risk data to personalize interactions.
- Train support agents to recognize churn risk indicators and respond with empathy and appropriate solutions.





## 6.4. Strategic Implementation Plan

| Action                           | Timeline    | Responsibility           | Key Metrics                           |
|----------------------------------|-------------|--------------------------|---------------------------------------|
| Launch gender-specific campaigns | 3 months    | Marketing Team           | Churn rate by gender                  |
| Develop age-based loyalty offers | 3-6 months  | Customer Experience Team | Retention rate among older customers  |
| Implement region-specific plans  | 6 months    | Regional Managers        | Churn rate in Germany                 |
| Product bundling offers          | 3 months    | Product Management Team  | Number of multi-product customers     |
| Reactivation campaigns           | 3 months    | Customer Engagement Team | Activity rate among dormant accounts  |
| Neural network deployment        | 6-12 months | Data Science Team        | Accuracy and ROC of churn predictions |

## 7. Conclusion

To effectively reduce churn, RBC must adopt a **comprehensive, data-driven approach** that focuses on personalized solutions, targeted engagement, and strategic testing. By leveraging insights from predictive models, RBC can:

1. **Engage Older Customers:** Offer diverse product sets that cater to the unique needs of older demographics, such as bundled mortgages, investment services, and retirement planning tools.
2. **Support Female Customers:** Develop gender-specific strategies that foster trust, financial empowerment, and personalized service for female clients.
3. **Reactivate Inactive Customers:** Launch targeted campaigns to incentivize dormant customers to re-engage with RBC services.
4. **Test and Refine Strategies:** Implement pilot programs and A/B testing to ensure that initiatives are effective and resource efficient.

By addressing the specific needs of different customer segments and continuously refining strategies based on data, RBC can improve customer satisfaction, enhance loyalty, and drive long-term business growth. This proactive approach not only reduces churn but also strengthens RBC's position as a customer-centric financial institution.



## References

Bailyn, E. (2024, August 1). Average customer acquisition cost (CAC) in banking. First Page Sage. <https://firstpagesage.com/seo-blog/average-customer-acquisition-cost-cac-in-banking/>