# DR KHANH
**Writing Service**

# BROCHURE

## The World - Class Assignment Service

*That you deserve*

**CONTACT US**

- DrKhanhAssignmentService
- www.drkhanh.edu.vn
- (+84) 939 070 595 hoặc (+84) 348 308 628

>>

# 1. Introduction

This report delineates a comprehensive data processing plan for analyzing customer churn within a financial institution. The primary objective is to predict which existing customers are most likely to discontinue their services, thereby facilitating the development of a targeted, evidence-based marketing strategy to enhance customer retention. The plan elucidates the methodological approach for processing and analyzing the provided dataset using IBM SPSS Modeler, in alignment with contemporary data mining practices in customer relationship management (Ngai et al., 2009).

Customer churn, defined as the propensity of customers to cease doing business with a company, represents a significant challenge in the financial services sector. As Keramati et al. (2016) note, the cost of acquiring new customers often substantially exceeds that of retaining existing ones, underscoring the importance of effective churn prediction and prevention strategies. This data processing plan aims to leverage advanced analytics techniques to identify at-risk customers and provide actionable insights for retention efforts.

# 2. Overview of the Data

The dataset comprises information on bank customers, encompassing eleven variables: CreditScore, Geography, Gender, Age, Tenure, Balance, NumOfProducts, HasCrCard, IsActiveMember, EstimatedSalary, and Exited. This multivariate dataset exemplifies the complexity of customer-related data in the banking sector, as noted by Larivière and Van den Poel (2005). The "Exited" variable serves as the dependent variable for churn prediction, with binary classification (1 indicating a churned customer, 0 indicating a retained customer).

The diversity of variables in this dataset allows for a multifaceted analysis of factors influencing customer churn. Demographic variables such as Geography, Gender, and Age provide insights into customer segments with higher churn propensity. Financial indicators like CreditScore, Balance, and EstimatedSalary offer a view into the economic aspects of customer behavior. Engagement metrics such as Tenure, NumOfProducts, HasCrCard, and IsActiveMember reflect the depth and breadth of a customer's relationship with the bank.

Understanding the interplay between these variables is crucial for developing a nuanced model of customer churn. As Tsai and Lu (2009) argue, comprehensive churn prediction models should consider a wide range of customer attributes to capture the complex dynamics of customer loyalty and attrition.

## 3. Data Pre-processing

Data pre-processing is a critical step in ensuring the quality and reliability of subsequent analyses (García et al., 2015). The following pre-processing steps will be undertaken:

- Data Cleaning: This phase will involve checking for missing values and outliers. Missing data, if any, will be addressed using multiple imputation techniques as recommended by Sterne et al. (2009). Outlier detection will be performed using the Interquartile Range (IQR) method, with treatment decisions based on the nature of the outlier and its potential impact on the analysis (Hodge and Austin, 2004).
- Encoding Categorical Variables: Categorical variables such as Geography and Gender will be encoded using one-hot encoding, a method widely used in machine learning for handling nominal categorical features (Potdar et al., 2017).
- Normalization: Numerical variables will be normalized using min-max scaling to ensure they are on a common scale, which is crucial for certain classification algorithms (Al Shalabi et al., 2006).

## 4. Defining Variables

In accordance with the provided guidelines, the dependent variable will be:

- Exited (0 = retained, 1 = churned)

All other variables in the dataset will be considered as potential predictors, aligning with the comprehensive approach to churn prediction advocated by Tsai and Lu (2009).

## 5. Calculating Churn Rate

The overall churn rate will be calculated using the formula:

$$Churn\ Rate\ =\ (Number\ of\ Churned\ Customers\ /\ Total\ Number\ of\ Customers)\ *\ 100$$

This metric provides a fundamental understanding of customer attrition and serves as a baseline for evaluating the effectiveness of retention strategies (Neslin et al., 2006).

# 6. Classification Model Application

We will employ five classification techniques available in SPSS Modeler, each with its unique strengths in predictive modeling:

- C5.0: An evolution of the C4.5 algorithm, C5.0 constructs decision trees or rule sets. It is known for its efficiency and handling of both continuous and categorical variables (Quinlan, 2014). C5.0 offers advantages in interpretability, making it particularly useful for understanding the decision paths leading to churn predictions.

- Tree-AS node: This algorithm creates decision trees using exhaustive CHAID (Chi-squared Automatic Interaction Detector) analysis, which is particularly effective in exploring relationships between variables (Kass, 1980). The Tree-AS node is adept at handling large datasets and can automatically detect complex interactions between predictors.

- Bayesian Network: This probabilistic graphical model represents variables and their conditional dependencies, offering interpretable results and handling of uncertainty (Friedman et al., 1997). Bayesian networks are particularly useful in situations where understanding the relationships between variables is as important as the predictions themselves.

- CHAID node: The Chi-squared Automatic Interaction Detector creates a decision tree using chi-square tests to identify optimal splits, excelling in situations with categorical predictors (Kass, 1980). CHAID can reveal multi-way splits in the data, potentially uncovering complex patterns in customer behavior.

- Neural Network: This model, inspired by biological neural networks, is capable of capturing complex non-linear relationships between variables (Zhang, 2000). While potentially less interpretable than other models, neural networks can often achieve high predictive accuracy, especially with large datasets.

For each model, we will:

1. Partition the data into training (70%) and testing (30%) sets. This split allows for model development on a substantial portion of the data while reserving a significant amount for unbiased evaluation.

2. Train the model on the training data. During this phase, we will pay careful attention to model parameters. For decision tree models (C5.0, Tree-AS, CHAID), we will experiment with different tree depths and splitting criteria. For the neural network, we will test various architectures and activation functions.

3. Generate predictions on the test set. This step provides an unbiased assessment of model performance on unseen data.

4. Evaluate performance using metrics including accuracy, precision, recall, F1-score, and ROC AUC. While accuracy provides an overall measure of correct predictions, precision and recall offer insights into the model's performance on positive (churned) cases. The F1-score balances precision and recall, while ROC AUC assesses the model's ability to distinguish between classes across various threshold settings.

To ensure robust performance estimates, k-fold cross-validation (k=5) will be employed, as recommended by Kohavi (1995). This technique provides a more reliable estimate of model performance by using multiple train-test splits.

# 7. Predicting Churn in New Dataset

Upon selection of the best-performing model based on the aforementioned metrics, we will:

1. Apply identical pre-processing steps to the new dataset. This ensures consistency between the training data and the new data on which predictions will be made.

2. Utilize the trained model to predict churn probabilities for each customer. Rather than simply classifying customers as 'churn' or 'not churn', we will use probability outputs to provide a more nuanced view of churn risk.

3. Classify customers into risk categories based on their predicted churn probabilities. We may define categories such as 'Low Risk' (0-30% churn probability), 'Medium Risk' (31-70% churn probability), and 'High Risk' (71-100% churn probability). These categories can guide the prioritization of retention efforts.

This approach aligns with best practices in predictive analytics for customer churn, as outlined by Verbeke et al. (2012). By providing probabilistic outputs and risk categories, we offer actionable insights that can be directly incorporated into customer retention strategies.

# 8. Visualizing Results

Data visualization plays a crucial role in communicating complex analytical results (Friendly, 2008). We will create the following visualizations:

1. Churn rate by customer segment (bar charts): This will visually represent the churn rates across different customer segments, helping to identify high-risk groups.

2. Feature importance plot for the best model: This visualization will highlight which variables have the strongest influence on churn predictions, guiding feature selection and interpretation.

3. ROC curve comparing model performance: By plotting the true positive rate against the false positive rate for various threshold settings, we can visually compare the performance of our different models.

4. Decision tree visualization (if applicable): For models like C5.0 or CHAID, a visual representation of the decision tree can provide intuitive insights into the prediction process.

These visualizations will be created using SPSS Modeler and supplemented with Python's matplotlib and seaborn libraries where necessary. Each visualization will be accompanied by a brief interpretation, ensuring that the insights are accessible to both technical and non-technical stakeholders.

# 9. Conclusion

This data processing plan outlines a rigorous approach to analyzing and predicting customer churn. By adhering to this methodology, we will identify key factors contributing to churn and segment customers based on their propensity to leave. The insights derived from this analysis will be instrumental in developing a targeted, evidence-based marketing strategy to improve customer retention, a critical factor in maintaining competitive advantage in the financial services sector (Keramati et al., 2016).

The multi-model approach allows us to leverage the strengths of different classification techniques, potentially uncovering complex patterns in customer behavior. By combining these advanced analytics with clear visualizations and interpretations, we aim to provide actionable insights that can directly inform retention strategies.

It is important to note that while this plan provides a robust framework for churn prediction, the effectiveness of any resulting strategies will depend on their implementation and the broader business context. Regular model updates and performance monitoring will be crucial to ensure ongoing relevance and accuracy in churn predictions.

# References

Al Shalabi, L., Shaaban, Z., & Kasasbeh, B. (2006). Data mining: A preprocessing engine. Journal of Computer Science, 2(9), 735-739.

Friedman, N., Geiger, D., & Goldszmidt, M. (1997). Bayesian network classifiers. Machine Learning, 29(2), 131-163.

Friendly, M. (2008). A brief history of data visualization. In Handbook of data visualization (pp. 15-56). Springer, Berlin, Heidelberg.

García, S., Luengo, J., & Herrera, F. (2015). Data preprocessing in data mining. Springer.

Hodge, V., & Austin, J. (2004). A survey of outlier detection methodologies. Artificial Intelligence Review, 22(2), 85-126.

Kass, G. V. (1980). An exploratory technique for investigating large quantities of categorical data. Applied Statistics, 29(2), 119-127.

Keramati, A., Jafari-Marandi, R., Aliannejadi, M., Ahmadian, I., Mozaffari, M., & Abbasi, U. (2016). Improved churn prediction in telecommunication industry using data mining techniques. Applied Soft Computing, 24, 994-1012.

Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model selection. In Ijcai (Vol. 14, No. 2, pp. 1137-1145).

Larivière, B., & Van den Poel, D. (2005). Predicting customer retention and profitability by using random forests and regression forests techniques. Expert Systems with Applications, 29(2), 472-484.

Neslin, S. A., Gupta, S., Kamakura, W., Lu, J., & Mason, C. H. (2006). Defection detection: Measuring and understanding the predictive accuracy of customer churn models. Journal of Marketing Research, 43(2), 204-211.

Ngai, E. W., Xiu, L., & Chau, D. C. (2009). Application of data mining techniques in customer relationship management: A literature review and classification. Expert Systems with Applications, 36(2), 2592-2602.

Potdar, K., Pardawala, T. S., & Pai, C. D. (2017). A comparative study of categorical variable encoding techniques for neural network classifiers. International Journal of Computer Applications, 175(4), 7-9.

Sterne, J. A., White, I. R., Carlin, J. B., Spratt, M., Royston, P., Kenward, M. G., ... & Carpenter, J. R. (2009). Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. BMJ, 338, b2393.

Tsai, C. F., & Lu, Y. H. (2009). Customer churn prediction by hybrid neural networks. Expert Systems with Applications, 36(10), 12547-12553.

Verbeke, W., Dejaeger, K., Martens, D., Hur, J., & Baesens, B. (2012). New insights into churn prediction in the telecommunication sector: A profit driven data mining approach. European Journal of Operational Research, 218(1), 211-229.

Zhang, G. P. (2000). Neural networks for classification: a survey. IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews), 30(4), 451-462.