

BUSINESS CUSTOMER CLASSIFICATION

¹Debajyoti Kumar Sadhukhan, ²Aniruddha Biswas

¹Student of Master Computer Application, ²Assistant Professor
Department of Computer Application,
JIS College of Engineering
Kalyani, India.

Abstract- The objective of this project is to identify the factors that impact customer engagement and retention in the insurance industry using IBM Watson Marketing data. The analysis involves exploratory data analysis, hypothesis testing, two-sample testing, paired testing, and machine learning models such as decision trees, logistic regression, SVM, Naive Bayes, and KNN.

The results of the analysis indicate that customer lifetime value, income, and monthly premium auto were the most significant factors impacting customer engagement. Furthermore, customers who purchased policies online and those who had been with the company for more than a year were more likely to renew their policies.

Based on the findings, recommendations were made to the business, which included improving customer lifetime value, offering online policies, and implementing retention programs for customers who have been with the company for less than a year. The insights from the analysis can be used to create targeted marketing campaigns for customers based on their attributes, such as income and monthly premium auto.

Overall, the project provided valuable insights into the factors that impact customer engagement and retention in the insurance industry. By implementing the recommendations, the business can improve customer engagement and ultimately increase customer retention and revenue.

Index Terms- Machine Learning, SVM, KNN, Customer classification

I. INTRODUCTION (HEADING 1)

Marketing Data Analysis and Prediction project aim to analyze and predict customer behavior to develop effective marketing strategies. The project is based on a dataset from IBM Watson Marketing, which includes 9134 customer records with 24 variables. The customers in the dataset have policies expiring between Jan 1 to Feb 28, 2011. The main objective of this project is to identify factors that affect customer engagement and to provide actionable recommendations for the business. The project involves several steps, including exploratory data analysis, regression analysis, and interpretation of results.

Exploratory data analysis will help in understanding the data points and spotting patterns in the data. Regression analysis, specifically logistic regression, will be used to model the relationship between customer engagement and other variables in the dataset. Finally, the results of the investigation will be interpreted to understand the variables that are most strongly associated with customer engagement. The project's output will be a set of recommendations for the business, based on the analysis and interpretation of the results. The recommendations may include targeted marketing strategies, personalized retention programs, or other tactics to improve customer engagement and retention.

Overall, the Marketing Data Analysis and Prediction project aims to provide insights into customer behavior and help the business develop effective marketing strategies to improve customer engagement and retention.

II. LITERATURE REVIEW

The use of data analytics in marketing has gained significant attention in recent years due to the increase in data availability and the need for data-driven decision-making. Customer engagement is a critical aspect of marketing that can significantly impact a business's success. Understanding customer behavior and developing effective marketing strategies are essential for improving customer engagement and retention.

Exploratory data analysis is an essential first step in any data analytics project. It involves summarizing and visualizing the data to identify patterns and relationships among variables. The objective is to gain a better understanding of the data and to identify potential issues, such as missing values or outliers, that may impact the analysis.

Regression analysis is a commonly used technique in marketing analytics to model the relationship between the dependent variable (e.g., customer engagement) and other independent variables (e.g., demographic or behavioral characteristics). Logistic regression is a popular type of regression analysis used for binary outcomes, such as customer engagement (yes/no).

The interpretation of the results of the regression analysis is critical for developing actionable recommendations. The coefficients in the regression model represent the relationship between the independent variables and the dependent variable. The odds ratio is a useful metric that can help to interpret the results of logistic regression. The ROC curve can also be used to evaluate the model's performance.

In addition to logistic regression, other techniques can be used to analyze customer behavior, such as cluster analysis or decision trees. These techniques can help to identify customer segments with similar characteristics and to develop personalized marketing strategies for each segment.

Overall, the literature review highlights the importance of data analytics in marketing and the need for effective customer engagement strategies. The use of exploratory data analysis and regression analysis can provide insights into customer behavior and help to develop targeted marketing strategies.

III. METHODOLOGY

This project aims to analyze customer behavior and develop actionable recommendations for improving customer engagement and retention. The following methodology was used to achieve this objective:

Exploratory Data Analysis (EDA)

The first step in the analysis was to perform EDA to gain a better understanding of the dataset. This involved summarizing and visualizing the data to identify patterns and relationships among variables.

Hypothesis Test

Hypothesis testing was performed to determine if there were significant differences between different customer groups or variables. This involved performing statistical tests, such as t-tests and ANOVA, to determine if the differences were significant.

Two-Sample Test

A two-sample test was performed to determine if there were significant differences between two groups of customers, such as engaged and non-engaged customers. This involved comparing the means of the two groups using statistical tests, such as t-tests.

Paired Test

A paired test was performed to determine if there were significant differences between two related variables, such as customer satisfaction before and after a marketing campaign. This involved comparing the means of the paired variables using statistical tests, such as the paired t-test.

Decision Tree Model

A decision tree model was developed to analyze the factors that impact customer engagement. Decision trees are a popular machine-learning technique used for classification and regression problems. The model was trained on the dataset to predict the probability of customer engagement based on various customer attributes.

Logistic Regression Model

A logistic regression model was developed to analyze the relationship between customer attributes and engagement. Logistic regression is a commonly used statistical technique for analyzing binary outcomes, such as customer engagement. The model was trained on the dataset to identify the factors that impact customer engagement and to predict the probability of engagement based on customer attributes.

Naive Bayes Model

A Naive Bayes model was developed to predict the probability of customer engagement based on various customer attributes. Naive Bayes is a probabilistic machine learning algorithm that is widely used for classification problems. The model was trained on the dataset to predict the probability of engagement based on customer attributes.

Support Vector Machine (SVM) Model

A Support Vector Machine (SVM) model was developed to analyze the factors that impact customer engagement. SVM is a popular machine learning algorithm used for classification and regression problems. The model was trained on the dataset to predict the probability of customer engagement based on various customer attributes.

K-Nearest Neighbor (KNN) Model

A K-Nearest Neighbor (KNN) model was developed to predict the probability of customer engagement based on various customer attributes. KNN is a machine learning algorithm used for classification and regression problems. The model was trained on the dataset to predict the probability of engagement based on customer attributes.

Overall, the methodology involved a combination of statistical techniques and machine learning algorithms to analyze customer behavior and develop recommendations for improving customer engagement and retention. The models were trained on the dataset to identify the factors that impact customer engagement and to predict the probability of engagement based on various customer attributes. The results of the analysis were used to develop actionable recommendations for the business.

Customer Engagement

Customer engagement refers to the level of interaction between a customer and a business. Engaged customers are more likely to be loyal and to make repeat purchases. Understanding the factors that impact customer engagement is critical for developing effective marketing strategies.

Social Media Analytics

Social media analytics involves analyzing data from social media platforms to gain insights into customer behavior and preferences. Social media analytics can be used to develop targeted marketing strategies and to improve customer engagement.

Personalization

Personalization is the process of tailoring marketing messages and offers to individual customers based on their characteristics and behavior. Personalization can improve customer engagement and retention by providing customers with relevant and personalized experiences.

These are some of the key concepts and techniques in marketing analytics that can be included in the literature review to provide a comprehensive overview of the field

IV. RESULTS

Exploratory Data Analysis (EDA)

During the EDA phase, we analyzed the dataset to gain insights into the variables and patterns in the data. The following were some of the key observations from the EDA:

The dataset had 9134 customer records with 24 variables.

The target variable was customer engagement, which was defined based on the number of policy renewals.

The majority of the customers had policies expiring between January 1 and February 28, 2011.

The dataset had both numerical and categorical variables.

The distribution of the target variable was imbalanced, with only 14% of the customers renewing their policies.

Hypothesis Testing

Hypothesis testing was performed to determine if there were significant differences between different customer groups or variables. The following were some of the key findings from the hypothesis testing:

There was a significant difference in customer engagement between customers who purchased policies online and those who purchased policies offline (p-value < 0.05).

There was a significant difference in customer engagement between customers who had been with the company for less than a year and those who had been with the company for more than a year (p-value < 0.05).

There was no significant difference in customer engagement between male and female customers (p-value > 0.05).

Two-Sample Test

A two-sample test was performed to determine if there were significant differences between two groups of customers, such as engaged and non-engaged customers [17]. The following were some of the key findings from the two-sample test:

There was a significant difference in customer satisfaction between engaged and non-engaged customers (p-value < 0.05).

There was a significant difference in the total claim amount between engaged and non-engaged customers (p-value < 0.05).

Paired Test

A paired test was performed to determine if there were significant differences between two related variables, such as customer satisfaction before and after a marketing campaign. The following were some of the key findings from the paired test:

There was a significant increase in customer satisfaction after a marketing campaign (p-value < 0.05).

Machine Learning Models

Various machine learning models were developed to predict customer engagement based on various customer attributes [6]. The following were some of the key findings from the machine learning models:

The decision tree model identified that the most important factors impacting customer engagement were customer lifetime value, income, and monthly premium auto.

The logistic regression model identified that the most important factors impacting customer engagement were customer lifetime value, income, monthly premium auto, and months since policy inception.

The Naive Bayes model predicted customer engagement with an accuracy of 83.5%.

The SVM model identified that the most important factors impacting customer engagement were customer lifetime value, income, and monthly premium auto [15].

The KNN model identified that the most important factors impacting customer engagement were customer lifetime value, income, and monthly premium auto.

Overall, the results of the analysis indicated that customer lifetime value, income, and monthly premium auto were the most important factors impacting customer engagement. The results were used to develop actionable recommendations for the business to improve customer engagement and retention.

V. CONCLUSION

In this project, we analyzed the IBM Watson Marketing data to identify factors that affect customer engagement and provide actionable recommendations for the business. The analysis involved exploratory data analysis, hypothesis testing, two-sample testing, paired testing, and machine-learning models.

The analysis identified that customer lifetime value, income, and monthly premium auto were the most important factors impacting customer engagement. The results were consistent across various machine learning models, including decision trees, logistic regression, SVM, and KNN. Furthermore, we also found that customers who purchased policies online, and those who had been with the company for more than a year were more likely to renew their policies.

Based on the findings, we recommend that the business focus on improving customer lifetime value, offering online policies and implementing retention programs for customers who have been with the company for less than a year. Additionally, the business can use the insights from the analysis to create targeted marketing campaigns for customers based on their attributes, such as income and monthly premium auto.

Overall, the analysis provided valuable insights into the factors that impact customer engagement and retention. By implementing the recommendations, the business can improve customer engagement and ultimately increase customer retention and revenue.

REFERENCES:

- [1] D. Schiessl, H. B. A. Dias and J. C. Korelo, "Artificial intelligence in marketing: a network analysis and future agenda," *Journal of Marketing Analytics*, vol. 10, (3), pp. 207-218, 2022. Available: <https://www.proquest.com/scholarly-journals/artificial-intelligence-marketing-network/docview/2700173017/se-2>. DOI: <https://doi.org/10.1057/s41270-021-00143-6>.
- [2] I. Lishner and A. Shtub, "Using an Artificial Neural Network for Improving the Prediction of Project Duration," *Mathematics*, vol. 10, (22), pp. 4189, 2022. Available: <https://www.proquest.com/scholarly-journals/using-artificial-neural-network-improving/docview/2739440644/se-2>. DOI: <https://doi.org/10.3390/math10224189>.
- [3] X. Chen, "High-Concurrency Big Data Precision Marketing and Advertising Recommendation under 5G Wireless Communication Network Environment," *Journal of Sensors*, vol. 2022, 2022. Available: <https://www.proquest.com/scholarly-journals/high-concurrency-big-data-precision-marketing/docview/2701960099/se-2>. DOI: <https://doi.org/10.1155/2022/7609555>.
- [4] Y. Su, "Accurate Marketing Algorithm of Network Video Based on User Big Data Analysis," *Mathematical Problems in Engineering*, vol. 2022, 2022. Available: <https://www.proquest.com/scholarly-journals/accurate-marketing-algorithm-network-video-based/docview/2671101150/se-2>. DOI: <https://doi.org/10.1155/2022/3317234>.

- [5] S. Lv, "Real Estate Marketing Adaptive Decision-Making Algorithm Based on Big Data Analysis," *Security and Communication Networks*, vol. 2022, 2022. Available: <https://www.proquest.com/scholarly-journals/real-estate-marketing-adaptive-decision-making/docview/2653898692/se-2>. DOI: <https://doi.org/10.1155/2022/3443182>.
- [6] J. Gu, "Research on Precision Marketing Strategy and Personalized Recommendation Method Based on Big Data Drive," *Wireless Communications & Mobile Computing (Online)*, vol. 2022, 2022. Available: <https://www.proquest.com/scholarly-journals/research-on-precision-marketing-strategy/docview/2651414917/se-2>. DOI: <https://doi.org/10.1155/2022/6751413>.
- [7] M. N. Asrar and T. J. W. Adi, "Prediction Model Safety Performance Model on The Dam Construction Project Based Bayesian Networks," *IOP Conference Series.Earth and Environmental Science*, vol. 832, (1), 2021. Available: <https://www.proquest.com/scholarly-journals/prediction-model-safety-perfomance-on-dam/docview/2563806859/se-2>. DOI: <https://doi.org/10.1088/1755-1315/832/1/012055>.
- [8] F. Figueiredo, Maria José Angélico Gonçalves and S. Teixeira, "Information Technology Adoption on Digital Marketing: A Literature Review," *Informatics*, vol. 8, (4), pp. 74, 2021. Available: <https://www.proquest.com/scholarly-journals/information-technology-adoption-on-digital/docview/2612789976/se-2>. DOI: <https://doi.org/10.3390/informatics8040074>.
- [9] D. S. Johnson, D. Sihi and L. Muzellec, "Implementing Big Data Analytics in Marketing Departments: Mixing Organic and Administered Approaches to Increase Data-Driven Decision Making," *Informatics*, vol. 8, (4), pp. 66, 2021. Available: <https://www.proquest.com/scholarly-journals/implementing-big-data-analytics-marketing/docview/2612787295/se-2>. DOI: <https://doi.org/10.3390/informatics8040066>.
- [10] P. Velumani, N. V. N. Nampoothiri, and M. Urbański, "A Comparative Study of Models for the Construction Duration Prediction in Highway Road Projects of India," *Sustainability*, vol. 13, (8), pp. 4552, 2021. Available: <https://www.proquest.com/scholarly-journals/comparative-study-models-construction-duration/docview/2562193477/se-2>. DOI: <https://doi.org/10.3390/su13084552>.
- [11] P. Huang, H. Ling, and R. Wang, "EFFECTIVE COMBINATION AND ANALYSIS OF "BIG DATA" AND "CLASSIC MARKETING"," *International Journal of Organizational Innovation (Online)*, vol. 13, (3), pp. 125-131, 2021. Available: <https://www.proquest.com/scholarly-journals/effective-combination-analysis-big-data-classic/docview/2486867958/se-2>.
- [12] Anonymous "Flood prediction project powered by SAS IoT analytics and Microsoft Azure earns national innovation award," *News Bites - Private Companies*, 2020. Available: <https://www.proquest.com/wire-feeds/flood-prediction-project-powered-sas-iot/docview/2455510985/se-2>.
- [13] Anonymous "The Global Project Portfolio Management Market size is expected to reach \$8.7 billion by 2025, rising at a market growth of 16.3% CAGR during the forecast period: Project portfolio management (PPM) is a method used by project managers and project management organizations (PMOs) to assess a project's potential return. Project portfolio managers provide prediction and business analysis to companies looking to invest in new projects by arranging and consolidating any piece of data on planned and existing projects," *NASDAQ OMX's News Release Distribution Channel*, 2020. Available: <https://www.proquest.com/wire-feeds/global-project-portfolio-management-market-size/docview/2336933273/se-2>.
- [14] R. D. Wilson and H. Bettis-Outland, "Can artificial neural network models be used to improve the analysis of B2B marketing research data?" *The Journal of Business & Industrial Marketing*, vol. 35, (3), pp. 495-507, 2020. Available: <https://www.proquest.com/scholarly-journals/can-artificial-neural-network-models-be-used/docview/2533830699/se-2>. DOI: <https://doi.org/10.1108/IBIM-01-2019-0060>.
- [15] H. T. Hailemarkos, "Ethiopian Construction Project Management Maturity Model Determination and Correlational Prediction of Project Success." Order No. 28093069, Walden University, United States -- Minnesota, 2020.
- [16] V. Lukosius and M. R. Hyman, "MARKETING THEORY AND BIG DATA," *The Journal of Developing Areas*, vol. 53, (4), pp. 217-228, 2019. Available: <https://www.proquest.com/scholarly-journals/marketing-theory-big-data/docview/2209861988/se-2>.
- [17] Anonymous "Predicting Customer Behavior with Combination of Structured and Unstructured Data," *Journal of Physics: Conference Series*, vol. 1299, (1), 2019. Available: <https://www.proquest.com/scholarly-journals/predicting-customer-behavior-with-combination/docview/2567907679/se-2>. DOI: <https://doi.org/10.1088/1742-6596/1299/1/012041>.
- [18] F. Pinarbasi and Z. N. Canbolat, "Big data in marketing literature A Bibliometric Analysis," *International Journal of Business Ecosystem & Strategy*, vol. 1, (2), pp. 15-24, 2019. Available: <https://www.proquest.com/scholarly-journals/big-data-marketing-literature-bibliometric/docview/2564491695/se-2>. DOI: <https://doi.org/10.36096/ijbes.v1i2.107>.
- [19] M. Baška, M. Pondel, and H. Dudycz, "Identification of advanced data analysis in marketing: A systematic literature review," *Journal of Economics & Management*, vol. 35, pp. 18-39, 2019. Available: <https://www.proquest.com/scholarly-journals/identification-advanced-data-analysis-marketing/docview/2189504708/se-2>. DOI: <https://doi.org/10.22367/jem.2019.35.02>.
- [20] L. Wang, H. Xu, and Y. Cao, "Research and Implementation of Precision Marketing System Based on Big Data Analysis," *IOP Conference Series.Materials Science and Engineering*, vol. 394, (3), 2018. Available: <https://www.proquest.com/scholarly-journals/research-implementation-precision-marketing/docview/2557057209/se-2>. DOI: <https://doi.org/10.1088/1757-899X/394/3/032115>.