Joseph Arpaia, Victor Gonzalez, Matthew Susi, Nicholas Verendia, Ben Kosiek

Professor Lucy

Data Mining

30 April 2024

### *Part 4: Summary Paper and Discussion of Ethical Issues*

This paper will dig deep into customer demographics of the cruise line industry, relying highly on the data-based modeling techniques. This will adhere to the principles of the Cross-Industry Standard Process for Data Mining (CRISP-DM). Our overall objective for the upcoming project is the integration of advanced data analytics and customer modeling to create several groups based on the demographics, with the goal to implement personalized services and enhanced targeted marketing strategies.

The cruise industry, which is known for its special difficulties and possibilities, needs a holistic understanding of various draws, appeals and needs of the travelling class of people. The business industry can be very empowered if one can come up with and understand different patterns and classifications in the customer demographics. These strategies, on the other hand are not only conceived for the purpose of meeting customers' expectation but also for the purpose of enhancing operational efficiency and marketing effectiveness through which both these customers' satisfaction and business growth can be furthered.

Data preparation is a process that is both complicated and very important for our project. Cleaning the data and having good techniques to prepare it will create a strong basis for other parts of the analysis. The first stage of this process covered the rigorous editing of data quality

issues, and it was in this stage that we ensured that the missing values were well taken care of through the application of either imputation or deletion schemes in accordance with the appropriate data purity requirements. Besides, continuous variables were also converted to categorical variables to make data suitable for machine learning algorithms and guaranteed the right understanding and operation of data. Another manner of making the numerical data healthy was scaled grouping which minimizes a possible bias arise out of the different data ranges. Such detailed planning not only assured the above, but the preparatory work was very important for implementing the next exploratory and modeling process. Because of that our conclusions became valuable, and our recommendations could successfully be acted upon.

During the data exploration phase, we conducted comprehensive exploratory data analysis to examine the complexity of the customer demographic compositions, trends and behaviors of the cruise line industry. We made extensive use of models such as histograms in this process which helped a lot in revealing the key trends while picking up on the outliers. The impact of these visual tools also extended to the diverse characteristics and patterns in the customers and decomposed the data structure, which enabled the in-depth analysis of the objective. We were able to pinpoint these associations by performing data analysis on those graphical representations, and that helped our data modeling decisions further. This has been the sole step to evidently check whether the subsequent analyses have been based on data that are well understood and can be effectively clustered to develop the customer segmentation strategy.

During the data modeling phase of our project, we examined different machine learning models to understand the best way for us to perform the customer classification task based on cruise line. The tested models include Logistic Regression, Decision Tree Classifier, and K-Neighbors classifier. Logistic Regression yielded 80.8% training accuracy and 78.2% testing

accuracy. The Decision Tree Classifier yielded 98.5% training accuracy and 77.6% testing accuracy. K-Neighbors Classifier yielded 80.4% training accuracy and 67.5% testing accuracy. Where there are few models, the Decision Tree Classifier stands out as the best with perfect training accuracy of 98.5% and good testing accuracy of 77.6%. This model's success can be attributed to its ability to produce sophisticated decision surfaces which create a good fit in terms of the heterogenous nature of the customer demographic data that needs to be managed. The successful classification rate shows that the Tree Decision model was able to capture the data in effect, hence it was possible to come up with authentic predictions that can be used to promote targeted marketing of consumers and service personalization. This phase was important in determining the suitable analytical procedure, the building block towards useful insight that may guide business decisions in the industry of cruise lines.

Our project followed the CRISP-DM methodology, which structured our approach across all phases of data mining: from initial business understanding, through data exploration and preparation, to modeling and evaluation. Every step was performed with precision, allowing the data to be cleaned in the best way possible to help the models find the most effective ones. The CRISP-DM method helped the project in its daily running and ensured that its objectives were fair and linked to the strategic business goals, increasing our group project's success rate.

Classifying customers in the cruise line industry raises several ethical issues, notably privacy concerns and the risk of discrimination. Personalizing experiences based on demographic data can lead to the exclusion of individuals from certain offers or promotions, potentially resulting in unfair treatment. For example, those of lower classes or cheaper tickets, may have significantly limited experiences than those of first class. There is a possible chance of discrimination based on classifying customers. Furthermore, the use of collecting and analyzing

data to categorize them may possibly infringe on their right to privacy. It is crucial for companies to establish clear policies and ethical guidelines to manage and mitigate these risks, ensuring that customer data is used responsibly and inclusively.

This project has demonstrated the potential of data analytics in the cruise line industry, particularly in improving the understanding of customer demographics. We followed the guideline of CRISP-DM and this structure helped us to bring out the key findings that can be used for making decisions. These inferences offer us an excellent means of offering our clients personalized experiences, a factor which can enable us to improve our clients' happiness and business efficiency. The success of this project shows the impact that data-driven designs have in modifying business operations.