**Project Deliverable 4: Bank Marketing Dataset: EDA, Modeling, Clustering, and Pattern Mining**

Sagar Bhetwal, Umesh Dhakal, Nirajan Acharya, Suresh Ghimire, Enjal Chauhan

University of the Cumberlands

Department of Computer Science

Advanced Big Data and Data Mining(MSCS634)

Dr. Satish Penmatsa

February 15th, 2026

**Project Deliverable 4: Bank Marketing Dataset: EDA, Modeling, Clustering, and Pattern**

**Mining**

## Introduction

Our group chose a banking marketing dataset that looks at whether a customer decides to subscribe to a term deposit. We picked this dataset because it represents a real business problem and includes a large number of records with both numerical and categorical features. This made it a good fit for applying different data mining techniques covered in this course. The project was not just about building prediction models, but also about understanding the data, analyzing the results, and connecting what we found to real-world decision making.

## Why we Chose Banking Data

Our team selected the bank marketing dataset because it provides a realistic example of how companies use customer and campaign data to make decisions. The dataset contains a large number of records along with mixed data types (numeric and categorical), which made it useful for practicing data preprocessing, feature engineering, and different machine learning techniques. Since the outcome is whether a customer subscribed to a term deposit, it also supports clear evaluation and practical insights. One thing we liked about this dataset is that it has a mix of customer info and campaign info, so we can study both customer behavior and marketing strategy. Also, the target column is clear (subscribe or not), so it makes it easier to evaluate models and compare methods.
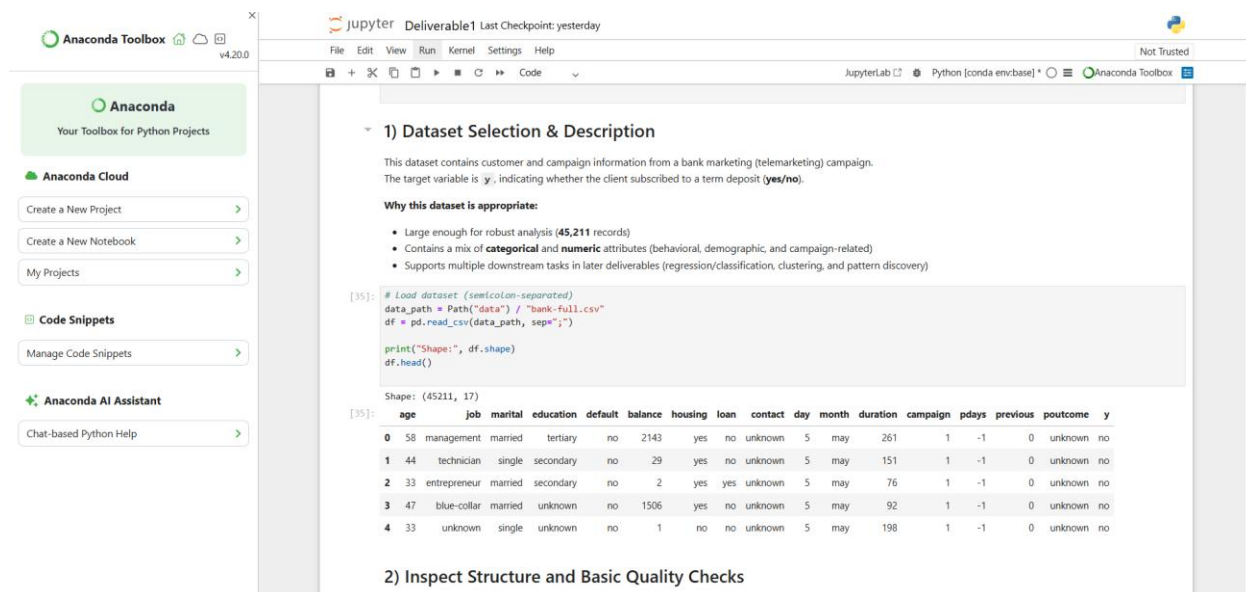
Figure 1. Bank marketing dataset

**Data Cleaning and Data Analysis**

During Deliverable 1, we focused on preparing the dataset so it could be used for modeling. We handled missing values carefully, especially in the pdays column. In this dataset, a missing pdays value means the customer was not contacted before. Instead of removing those rows, we created a flag variable to keep that information and filled the missing values in a way that would not affect the model. We also checked for duplicate records and looked at outliers in numerical features such as age, balance, duration, and campaign. Exploratory Data Analysis helped us better understand how the data is distributed, how features relate to each other, and how they connect to the target variable.

Some of the screenshots below show the target distribution, correlation heatmap, and numeric distributions.
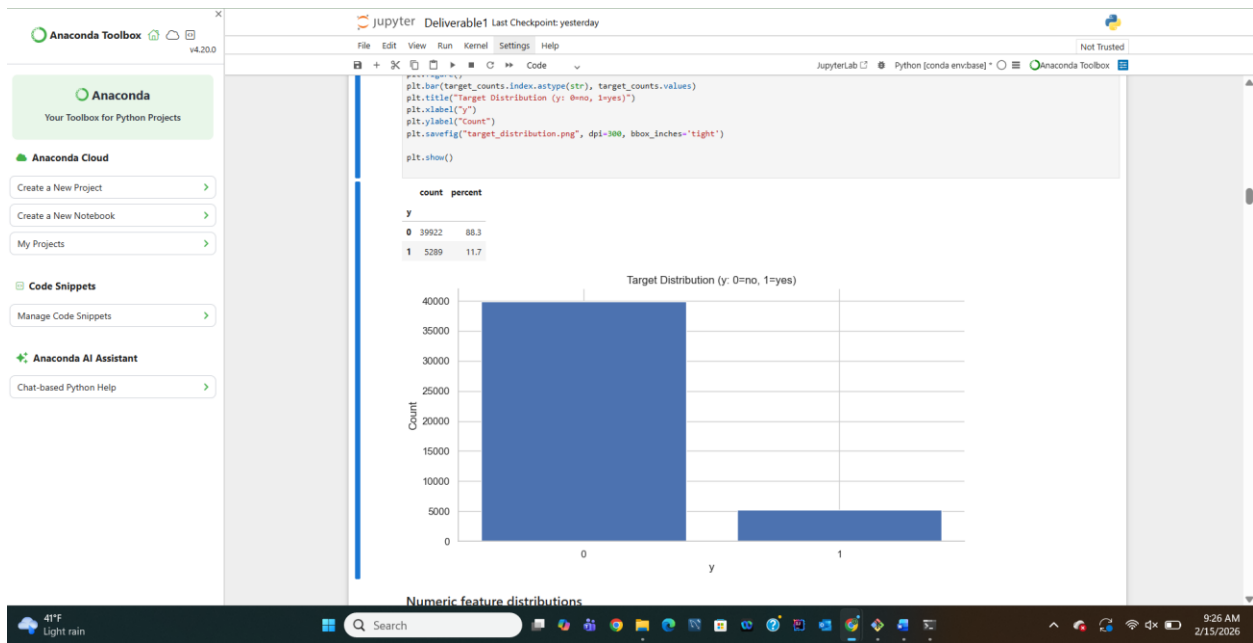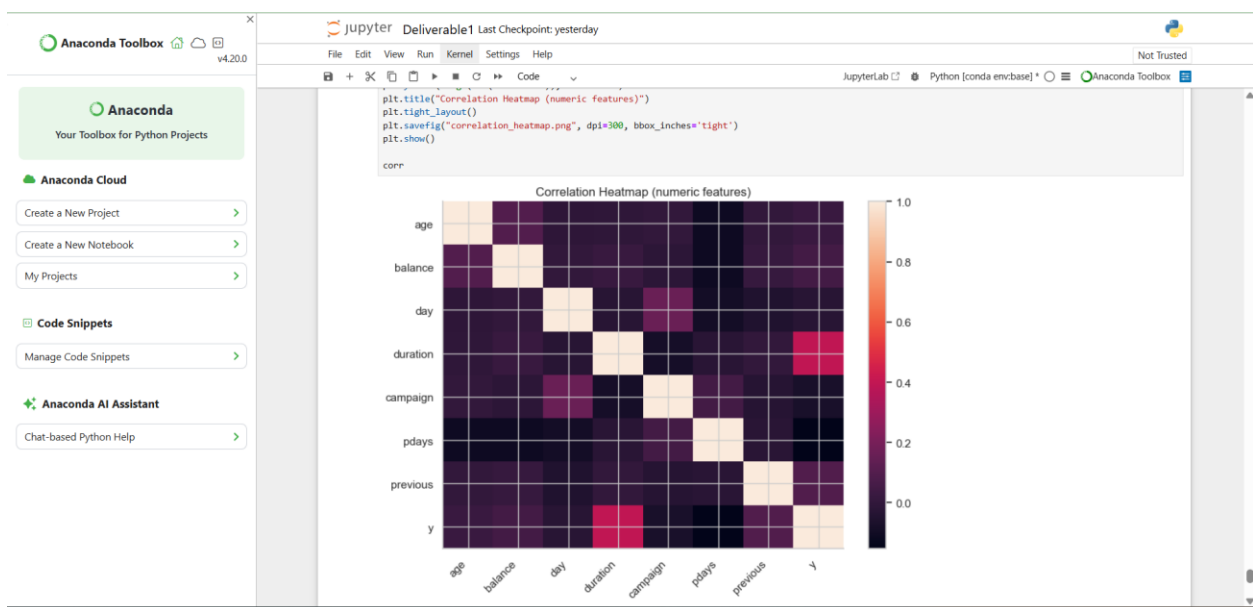
Figure 2. Subscription outcome distribution



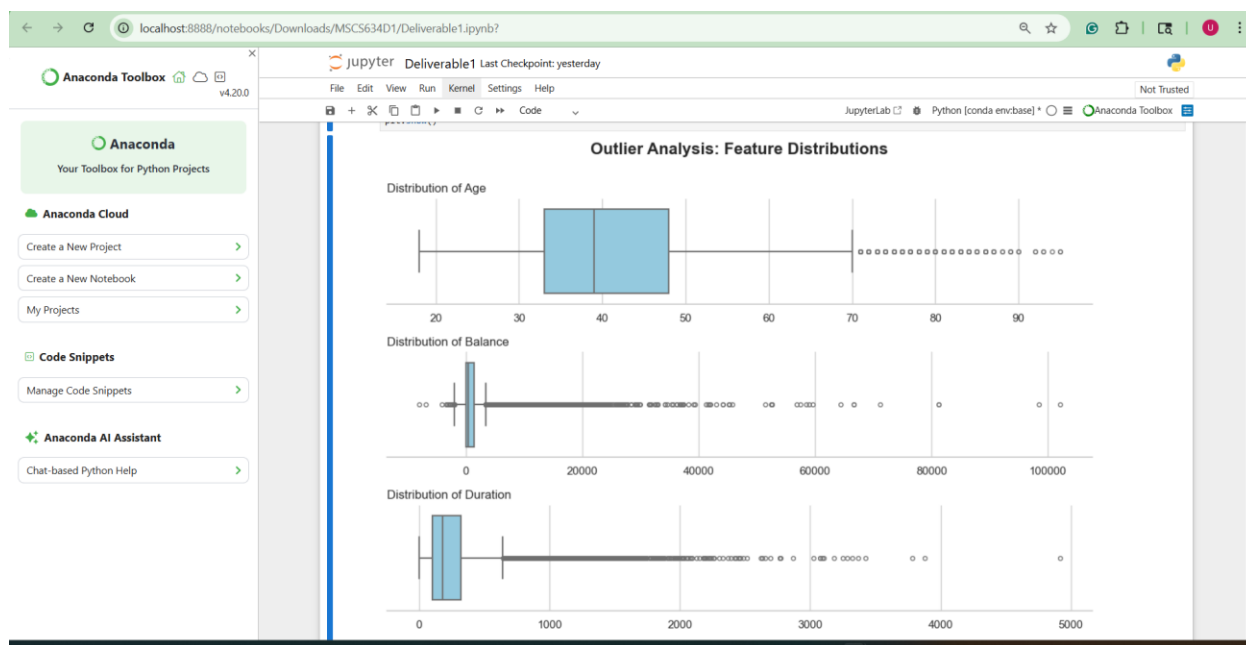Figure 3. Correlation heatmap of numeric features

Figure 4. Outlier and numeric feature distribution analysis

**Feature Engineering**

We did feature engineering to help improve the model performance. For example, we created age_squared so the model can capture a non-linear relationship between age and the target variable. We also made a few features related to previous campaign activity to better represent customer engagement and past contact history.

We also encoded categorical variables so they can work in machine learning models. These steps helped the models learn patterns from the data more clearly instead of ignoring important information.
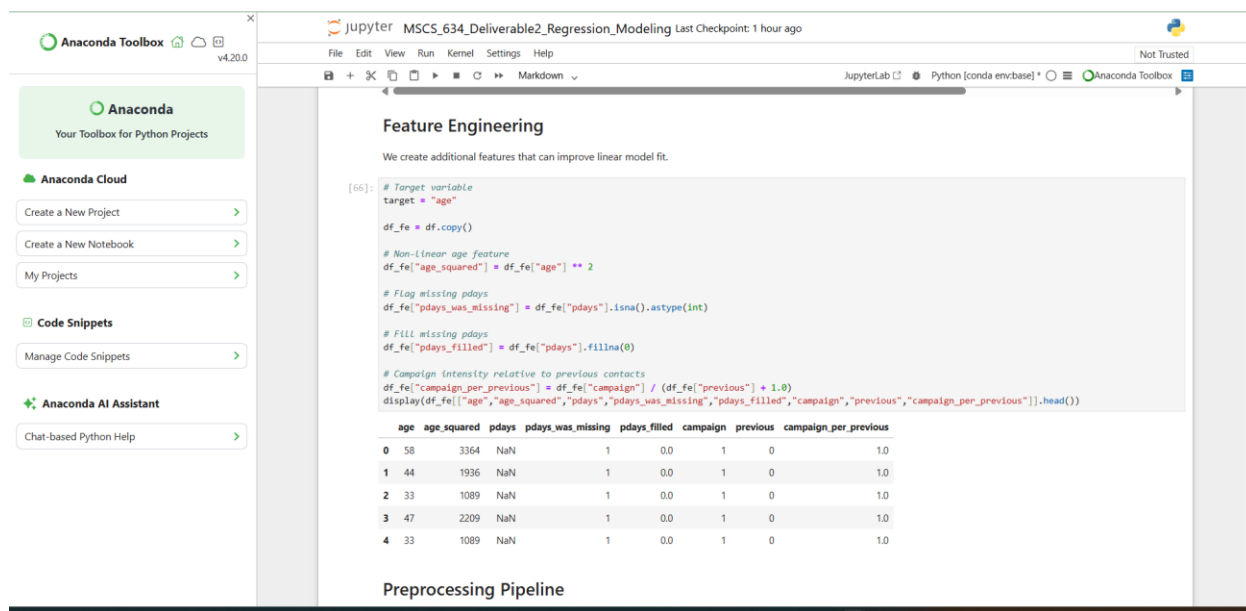
Figure 5. Feature engineering

## Regression Modeling and Evaluation

In Deliverable 2, our group created and compared multiple regression models to predict the subscription outcome. We started with baseline regression models and then improved them by adding engineered features. To measure how well each model performed, we used evaluation metrics such as R-squared, Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). We also used 5-fold cross-validation to test how consistent the models were across different splits of the data. This was helpful because it gave us a stronger idea of how the models might perform on new data and it reduced the chance of overfitting. The results table helped us see the difference between models in a simple way. Instead of picking a model from one split, cross-validation gave us a more consistent view of performance across multiple folds.
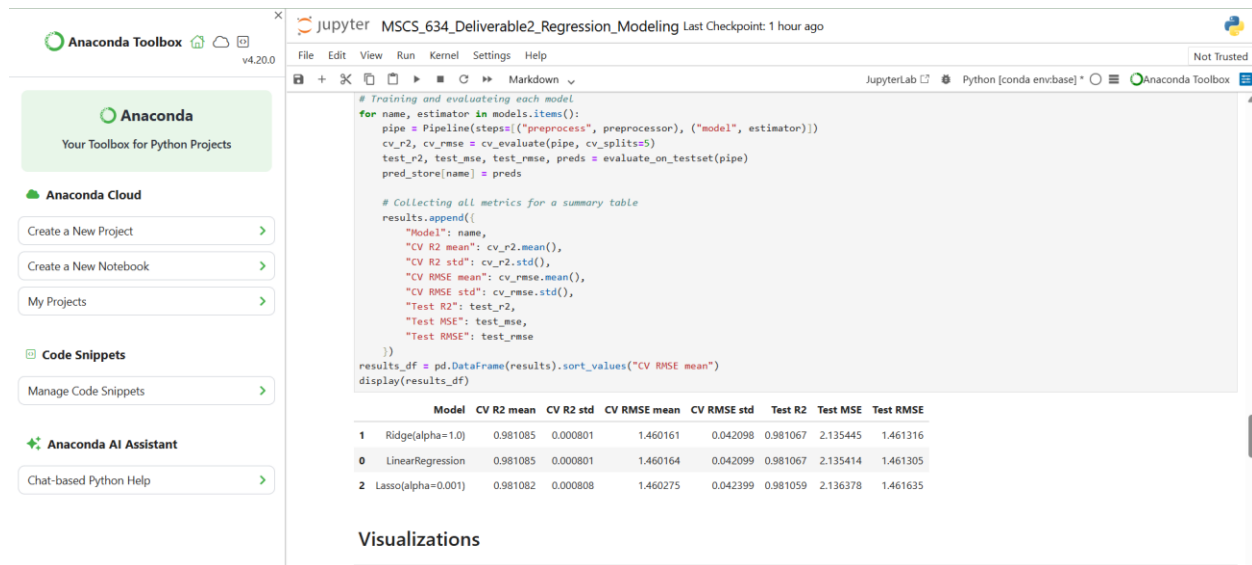
Figure 6. cross-validation/test results table.

**Classification, Clustering and Pattern Mining**

For Deliverable 3, our group worked on three main parts: classification, clustering, and pattern mining. First, we built multiple classification models such as Logistic Regression, Decision Trees, and Support Vector Machines (SVM) to predict whether a customer would subscribe. To compare the models, we used confusion matrices, ROC curves, accuracy, and F1 score. We also did hyperparameter tuning (mainly to improve the SVM model), which helped us get better results compared to the baseline settings. Next, we used K-Means clustering to group customers with similar characteristics. To choose the best number of clusters, we used the elbow method and silhouette scores, since they helped us see which cluster size made the most sense.

Finally, we applied association rule mining using the Apriori algorithm to find frequent patterns in the data. These rules helped us understand which combinations of customer features and campaign details were commonly linked with subscription outcomes, which can be useful for planning better marketing strategies. The confusion matrix and ROC curve made it easier to

see what kind of errors the model makes. This helped us understand whether the model was

missing too many "yes" cases or incorrectly predicting "yes" when it should be "no."
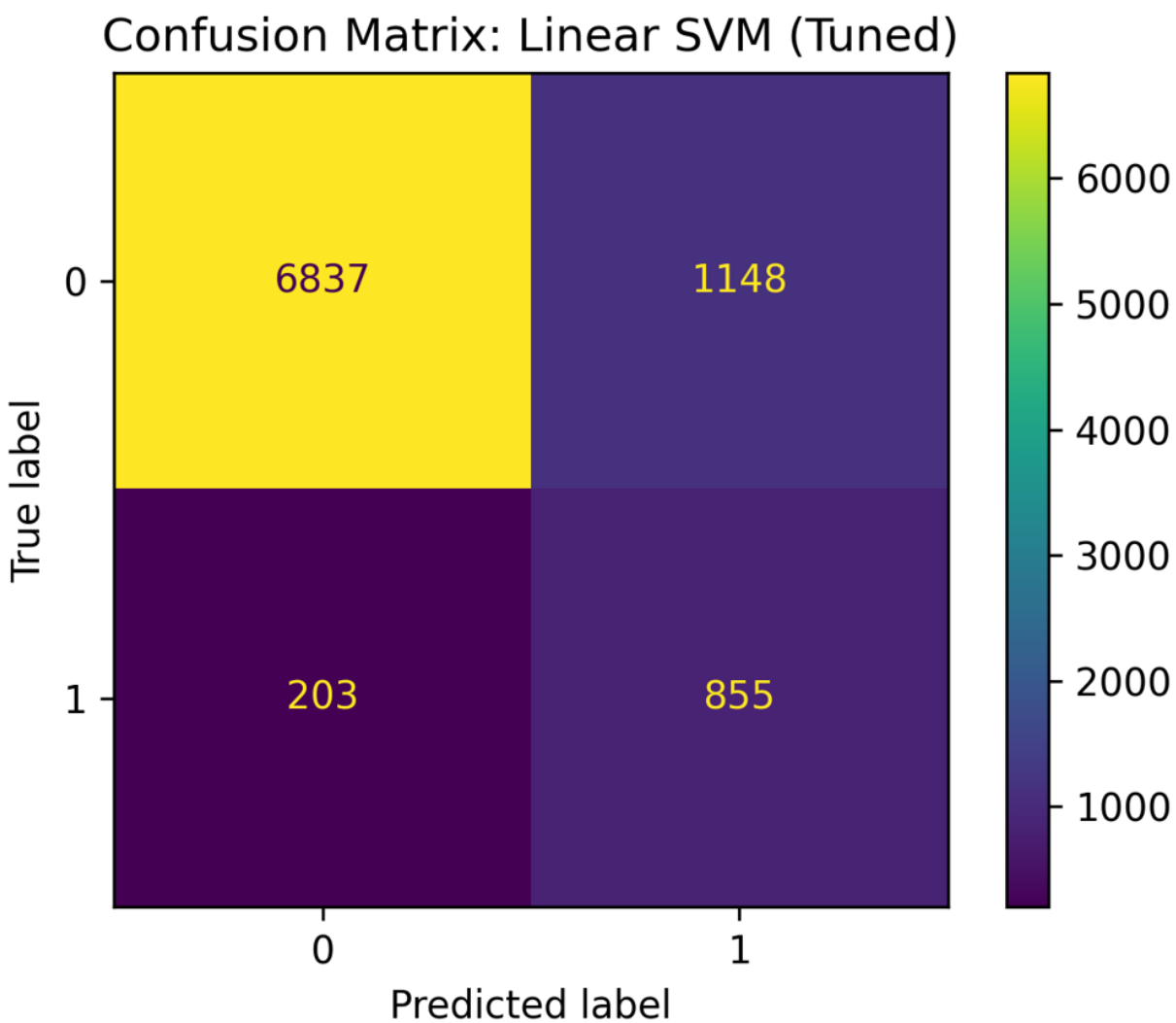


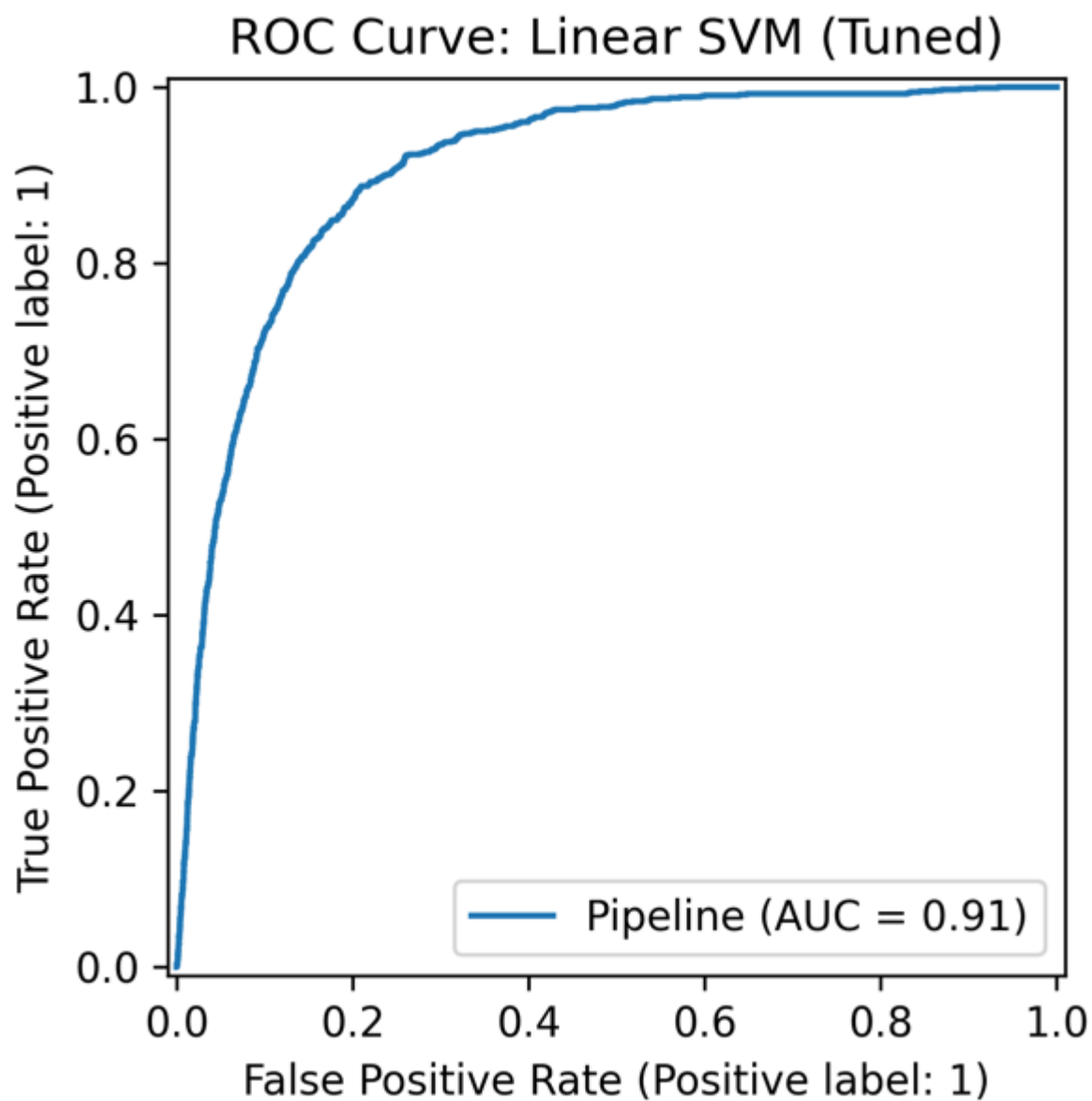Figure 7. Confusion matrix for tuned Linear SVM classifier

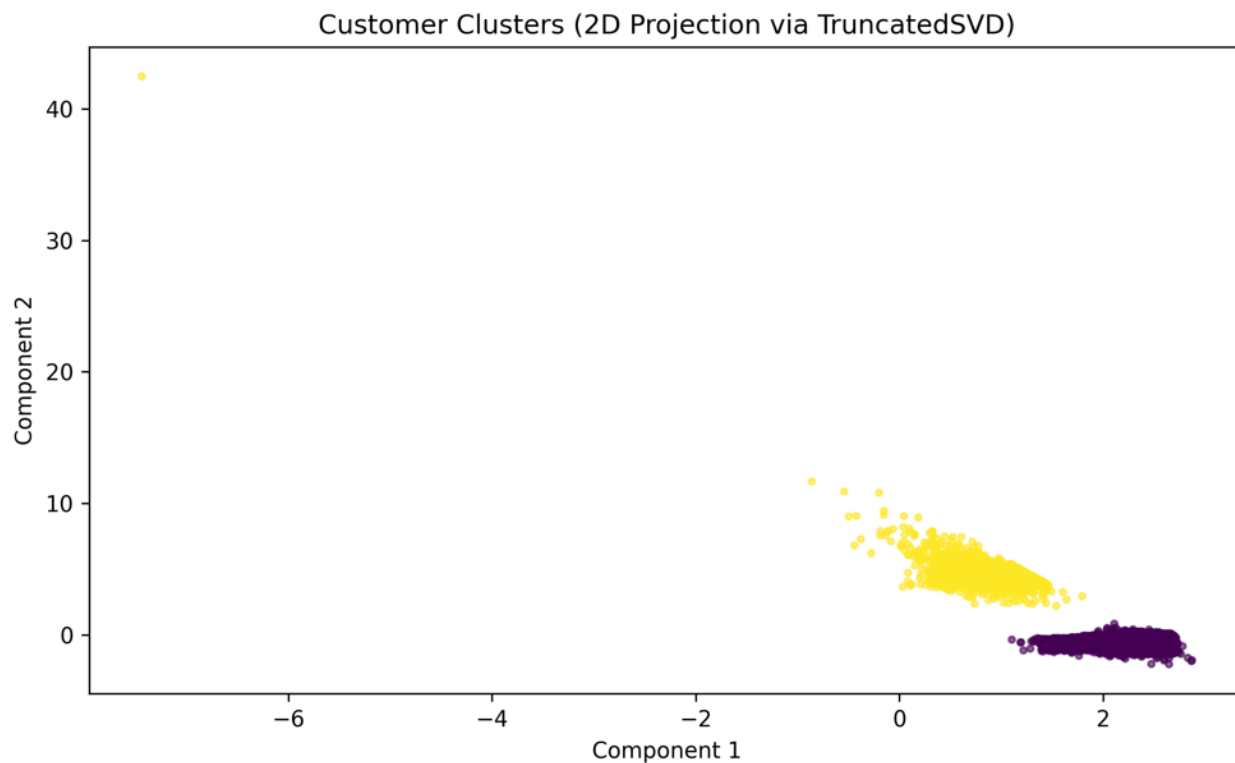Figure 8. ROC curve for tuned Linear SVM classifier
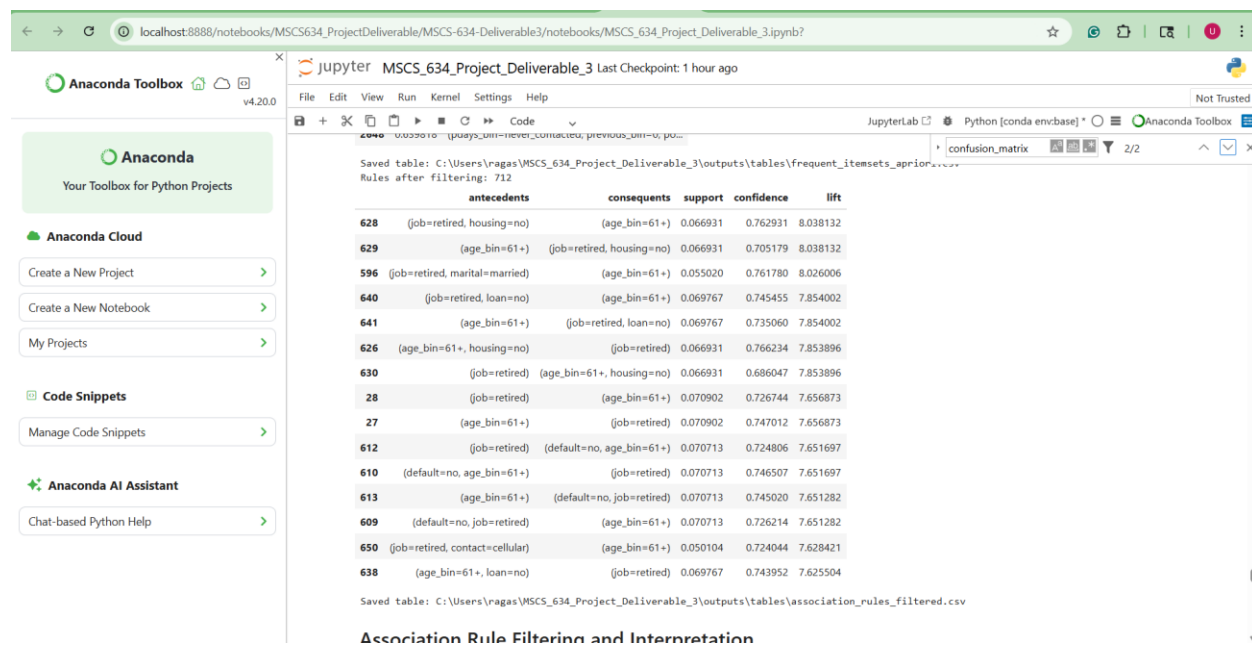
Figure 9. Customer clusters from K-Means



Figure 10. Top association rules from Apriori

**Key Insights We Found**

- Campaign features (especially duration and previous contact history) were some of the strongest signals for subscription.

- Classification models worked better than regression for this dataset, and results improved after hyperparameter tuning.

- K-Means clustering showed different customer groups, and some clusters had higher subscription rates than others.

- Association rules showed patterns where combinations of job type, education, and contact method appeared more often in successful subscriptions.

Overall, we noticed that campaign-related features like duration and previous contact history mattered a lot when predicting subscriptions. When we grouped customers using clustering, some segments were more likely to subscribe than others, which suggests that a targeted marketing strategy could work better than using the same approach for everyone. We also saw that classification models performed better than regression for this problem, especially after tuning the model settings. Finally, the association rules supported our findings by showing repeated combinations of customer attributes and contact method that showed up more often when subscriptions were successful.

**Ethical Considerations**

We kept ethical considerations in mind throughout this project. The dataset does not include personal identifying information like names, phone numbers, or addresses, which helps protect customer privacy. Even so, there can still be bias in the data because some features like age, job, and education might affect the outcome in a way that is not fully fair.

To focus on fairness, we paid attention to how the models behaved across different groups and we did not rely on features in a way that could lead to unfair or discriminatory decisions. Overall, our results should be used as a tool to support decisions, not as something that completely replaces human judgment.

**Future Work**

Based on our results, banks can make their marketing campaigns more effective by prioritizing customer groups that showed higher engagement and stronger subscription outcomes. Using past contact history and campaign response patterns to guide who to contact, when to contact, and how to contact them could improve the overall success rate. For future work, we could test more advanced approaches such as ensemble models like Random Forest or Gradient Boosting. It would also be useful to add time-based patterns such as seasonality, month trends and external factors like economic conditions to see if they improve prediction performance.

**Conclusion**

This project helped us go through the full data mining process from start to finish. We cleaned the data, handled missing values, and made sure everything was ready for modeling, then we did EDA to understand the distributions, relationships between features, and how different variables connect to the target outcome. After that, we built models and evaluated how well they worked, including both regression and classification so we could compare which approach fit this problem better. We also used K-Means clustering to group customers with similar characteristics, which helped us understand different customer segments, and we applied association rule mining to find common patterns that were often linked with subscription results.

# References

Dua, D., & Graff, C. (2019). UCI machine learning repository. University of California, Irvine, School of Information and Computer Sciences. https://archive.ics.uci.edu/ml/datasets/bank+marketing

Han, J., Kamber, M., & Pei, J. (2012). Data mining: Concepts and techniques (3rd ed.). Morgan Kaufmann.

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. Computing in Science & Engineering, 9(3), 90–95.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). An introduction to statistical learning: With applications in Python (2nd ed.). Springer.

Tan, P.-N., Steinbach, M., & Kumar, V. (2019). Introduction to data mining (2nd ed.). Pearson.