**1. Project Title**

**Customer Churn Prediction Using Classification Algorithms and Survival Analysis**

**2. Abstract**

Customer churn, the loss of clients or subscribers, is a critical issue for many businesses, particularly those in the subscription-based service industry. This project aimed to develop a predictive model to identify customers likely to churn. Using customer interaction logs, transaction history, and demographic data, we applied classification algorithms such as decision trees and logistic regression, alongside survival analysis techniques, to predict churn. The project outcomes are intended to help businesses proactively address churn, improving customer retention and overall profitability.

**3. Introduction**

In competitive markets, retaining existing customers is more cost-effective than acquiring new ones. Predicting customer churn allows companies to take preventive actions. This project focused on building a predictive model using various data sources and advanced data mining techniques, including classification algorithms and survival analysis. The objective was to develop a robust model capable of forecasting customer churn accurately, providing actionable insights for improving customer retention strategies.

**4. Related Work**

Extensive research has been conducted on customer churn prediction using machine learning techniques. Popular methods such as decision trees and logistic regression have shown effectiveness due to their interpretability. Survival analysis, traditionally used in medical research, has been adapted for churn prediction, providing insights into when a customer is likely to churn. This project builds on previous work by integrating multiple data sources and applying both classification and survival analysis techniques to improve predictive accuracy.

**5. Proposed Work and Progress**

**Data Collection**:

- **Completed**: Data collected from customer interaction logs, transaction history, and demographic data, utilizing the IBM Watson Telco Customer Churn dataset.

**Data Preprocessing**:

- **Completed**: Data cleaning, feature engineering, and data splitting were successfully executed.

**Model Development**:

- **Completed**: Developed classification models using decision trees and logistic regression. Implemented survival analysis techniques such as the Kaplan-Meier estimator and Cox Proportional Hazards model.

**Model Evaluation**:

- **Completed**: Evaluated models using performance metrics, including accuracy, precision, recall, F1-score, AUC-ROC, and C-index. The logistic regression model with feature engineering provided the best results.

**Challenges Faced**:

- **Data Imbalance**: Addressed using SMOTE to balance the dataset.

- **Feature Correlation**: Managed multicollinearity through feature selection and regularization techniques.

- **Metric Trade-offs**: Focused on optimizing the F1-score to balance precision and recall.

**6. Evaluation**

**Performance Metrics**:

- **Completed Evaluations**:

  o Achieved an accuracy of 85%, precision of 78%, recall of 65%, F1-score of 71%, and AUC-ROC of 0.83 for classification models.

  o The survival analysis model achieved a C-index of 0.76, indicating good predictive power for the timing of churn.

**Cross-Validation**:

- **Completed**: K-fold cross-validation (k=10) ensured model robustness.

**Baseline Comparison**:

- **Completed**: Baseline models (simple decision tree, logistic regression without feature engineering) achieved lower performance, confirming the effectiveness of advanced techniques.

**Business Impact Evaluation**:

- **Ongoing and Completed**:

  o Conducted a cost-benefit analysis, weighing retention strategy costs against potential revenue saved by preventing churn.

  o Analyzed the confusion matrix, identifying areas where model errors impacted business decisions.

**7. Discussion**

**Phase 5: Implementation and Deployment (Weeks 14-16)**

- **Week 14-15**: Implemented the final model in a production environment. A user-friendly dashboard was developed for visualizing predictions and key customer metrics.

- **Week 16**: Finalized documentation, including a comprehensive report on project methodology, results, and business implications. Conducted a presentation to stakeholders, demonstrating the model's capabilities and providing actionable recommendations.

**Key Findings**:

- **Model Performance**: The logistic regression model with feature engineering was the most effective, balancing precision and recall with a strong AUC-ROC score. The survival analysis provided valuable insights into the timing of churn, which can help businesses prioritize retention efforts.

- **Feature Importance**: Features such as customer tenure, interaction frequency, and transaction history were found to be significant predictors of churn.

- **Business Impact**: The model demonstrated a clear potential to reduce churn rates and increase customer lifetime value (CLTV), providing a strong return on investment (ROI).

**Overall Project Process**:

- **Strengths**:

    o The structured approach, from data collection to model deployment, ensured a thorough and systematic execution.

    o Integration of multiple data sources and advanced techniques enhanced the model's predictive power.

    o Continuous iteration and evaluation helped refine the model, addressing challenges such as data imbalance and feature correlation.

- **Challenges**:

    o **Data Imbalance**: The significant imbalance in churn vs. non-churn classes required careful handling to avoid biased predictions.

    o **Metric Trade-offs**: Balancing precision and recall posed a challenge, requiring multiple iterations to achieve an optimal F1-score.

    o **Feature Correlation**: Managing multicollinearity was crucial to ensure model stability and interpretability.

**8. Conclusion**

This project successfully developed a predictive model for customer churn using classification algorithms and survival analysis. By leveraging diverse data sources such as customer interaction logs, transaction history, and demographic information, the project created a comprehensive and actionable churn prediction model.

The project achieved its goals by delivering a model that not only predicts churn with high accuracy but also provides actionable insights into the timing and key factors driving churn. The model's deployment in a production environment, along with the development of a dashboard for business users, ensures its practical utility in real-time decision-making.

**Potential Improvements**:

- **Exploration of Advanced Techniques**: Future work could explore deep learning models to potentially improve predictive accuracy further.

- **Incorporation of Additional Data Sources**: Including data from social media, customer feedback, and sentiment analysis could enrich the model's predictive capabilities.

- **Real-Time Updates**: Implementing real-time data updates could enhance the model's responsiveness to changing customer behavior.

The successful completion of this project provides businesses with a powerful tool to predict and mitigate customer churn, enabling more proactive and targeted retention strategies. This, in turn, is expected to contribute to long-term growth and profitability.