# Telco Customer Churn Project

## Data Mining I Project

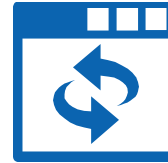Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

02.12.2024

# Overview

1. Introduction

2. Preprocessing

3. Finding common Customer Profiles

4. Machine Learning Model Analysis

5. Evaluation and Results

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# There goes a scenario:

One day, a customer called to complain, let's see how our model helps retain this customer…

# 1. Introduction

# Goals and Approach



**Telecom customer management grows together**

External: Turbulent and shifting industry

- Intensified competition
- Shifting consumer behavior

Internal: CRM capability challenges

- Action effectiveness
- Resource drain issues
- Customers acquisition challenge

**Core Issue**

Acquiring new customers tends to be **far more expensive** than retaining existing one

**Goals (What)**

Empower telecoms with predictive insights for **early action** and **stable customer retention**

**Actions (How)**

- Develop Model
- Apply ML
- Implement Strategies

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

6

# Dataset Structure

UNIVERSITÄT MANNHEIM

**91** features
*after splitting dataset

**67** Categorical

1 2 3
**24** continuous

a customer churn dataset

location

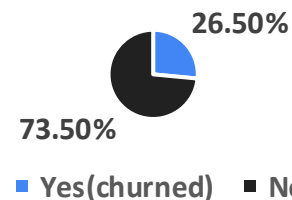population

**The dataset composition**

services

status

demographics

**Dataset overview**

- **7,000+**customer records
- The features spanning both categorical and numerical variables

**Dataset balance**

- The target variable, Churn, **is imbalanced**
- 26.5% labeled as Yes (churned) and 73.5% as No.

**26.50%**

**73.50%**

■ **Yes(churned)** ■ **No**

*This imbalance will be careful handling during model training, using techniques to ensure predictive fairness

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

# Hello, how can I help you?

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# Can you tell me your CustomerID?

## Telco Customer Churn Dashboard

1875-QIVME | Search

# What do we know about you?

## Telco Customer Churn Dashboard

| 1875-QIVME | Search |

Customer ID: 1875-QIVME

Churn Prediction: Likely to Churn

# 2. Preprocessing

# Preprocessing Pipeline



Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
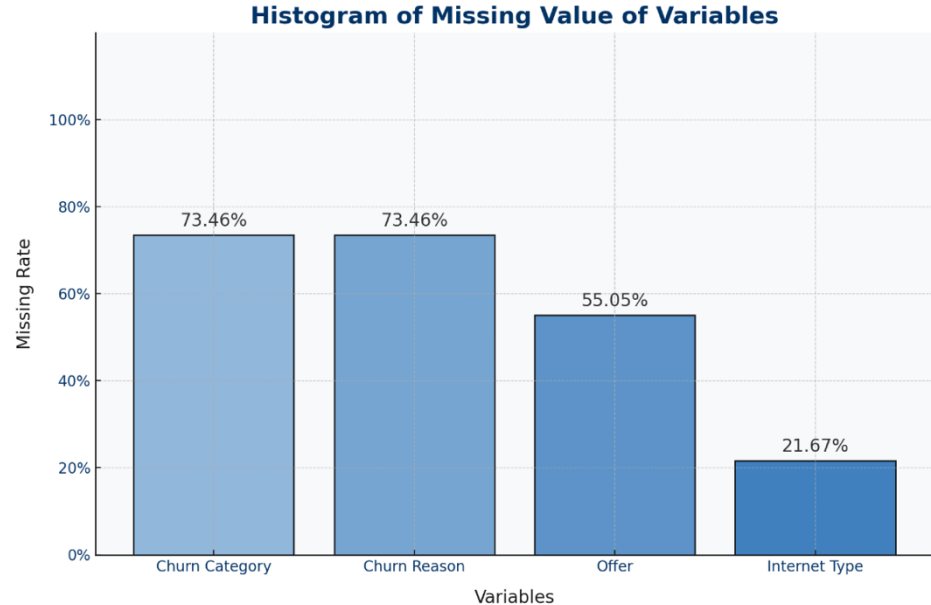
# Preprocessing

- **Handling Missing Values**
  - Significant Missing Values: substantial portion of the dataset affected, suggesting a strong relationship between the two columns.

  - Explanation: align with rows where the 'Churn' column is 'No', implying that these details were not recorded because they are not applicable to non-churning customers.



**Histogram of Missing Value of Variables**

# Preprocessing: Data Type Conversion

- **Removal of Unnecessary Columns**:
  - Removed columns with entirely unique entries (e.g., IDs) to retain only meaningful features.

- **Categorical Variable Encoding**:
  - **Label Encoding**: Applied to binary variables like "Gender" and "Senior Citizen."
  - **One-Hot Encoding**: Used for features with multiple categories

- **Post-Encoding Compatibility Check**:
  - Ensured all columns were numeric for compatibility with machine learning algorithms.

# Preprocessing: Feature Engineering

- **Key Approaches**
  - **Interaction Features**:
    - Combine variables (e.g., tenure vs. age) to reveal trends
  - **Aggregation**:
    - Summarize binary features (e.g., total services subscribed)
  - **Group-Based Features**:
    - Segment customers (e.g., revenue tiers by charges)
  - **Domain Transformations**:
    - Align features with business metrics (e.g., annualized charges, refund-to-charges ratio)

# Preprocessing

- **Data Splitting**
    - Split the data into training and test sets:
        - 80% (5,634 rows) for training
        - 20% (1,409 rows) for testing
- **Outlier Detection**
    - Outlier detection and removal are performed only on the training set
    - The data is scaled when applying each method
    - 3,885 rows remain in the training set after outlier removal

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# 3. Clustering

# Correlation Analysis

Long
Longitude:
1.0

Monthly Charges
Monthly_Charges_Scaled:
1.0

Tenure
Tenure_in_Years:
0.99982



Correlation Matrix for X_train

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
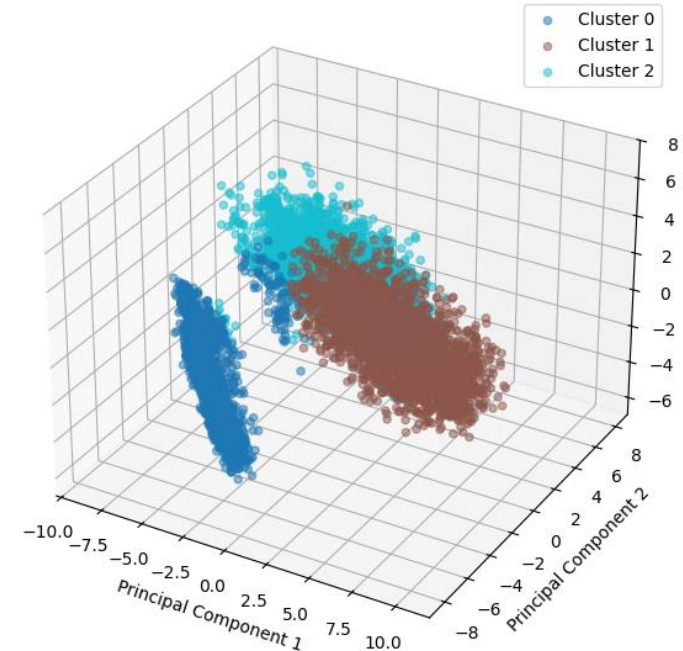02.12.2024

# Finding common Customer Profiles

Differentiators for the Churn Value – PCAs:

Loyality, Total Revenue, Total Charges

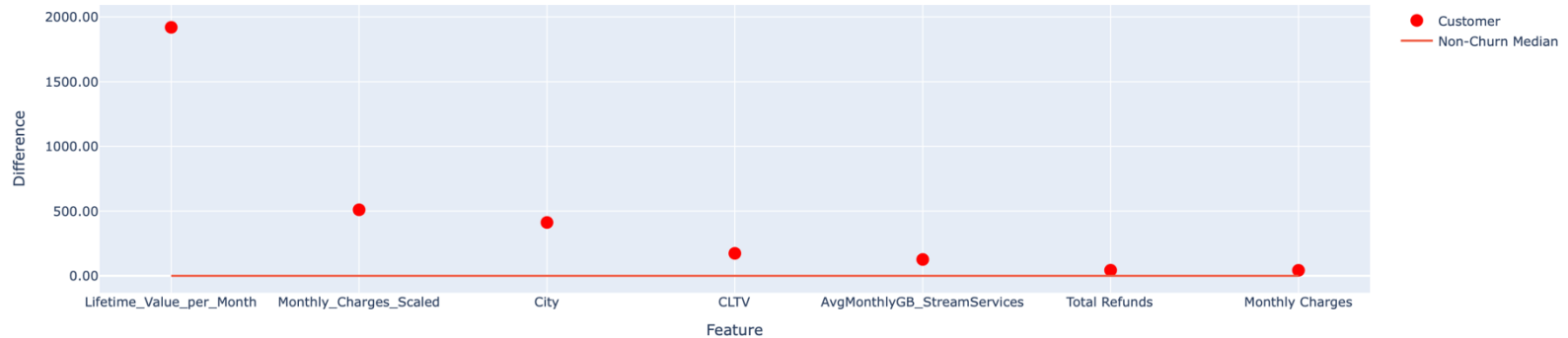| | |
|---|---|
| Cluster 0 (Churn 0.03) | Moderate revenue, low churn, **stable and satisfied** |
| Cluster 1 (Churn 0.08) | High revenue, low churn, **most valuable** |
| Cluster 2 (Churn 0.66) | Low revenue, high churn, **dissatisfied and at risk** |



3D PCA of Customer Data with K-Means Clusters

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# How we make data speak



Telco Customer Churn Dashboard

1875-QIVME | Search |

Customer ID: 1875-QIVME

Churn Prediction: Likely to Churn

Top 7 Feature Differences from Non-Churn Customers

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

20

# 4. Models Analysis

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

# Different Approaches

Naive Bayes

Logistic Regression

KNN

Nearest Centroid

Decision Tree

Random Forest

XG Boost

SVM

MLP

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

02.12.2024
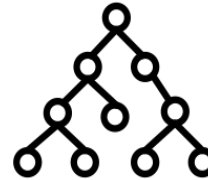
# 4 Main Types



## Probabilistic Models

- Gaussian Naive Bayes
- Multinomial Naive Bayes
- Logistic Regression

## Distance-Based Models

- K-Nearest Neighbors (KNN)
- Nearest Centroid

## Tree-Based Models

- Decision Tree
- Random Forest
- XG Boost

## Optimization and Kernel-Based Models

- Support Vector Machine (SVM)
- Multilayer Perceptron (MLP)

# Similar Procedures in 10 Models

I. Train Models

II. Hyperparameter Tuning

III. Cross-Validation

IV. Performance Metrics

# Probabilistic Models - Naive Bayes

## Gaussian Naive Bayes

I. **Feature Preparation**
   - Discretization of **continuous** features

II. Train Model

III. Hyperparameter Tuning
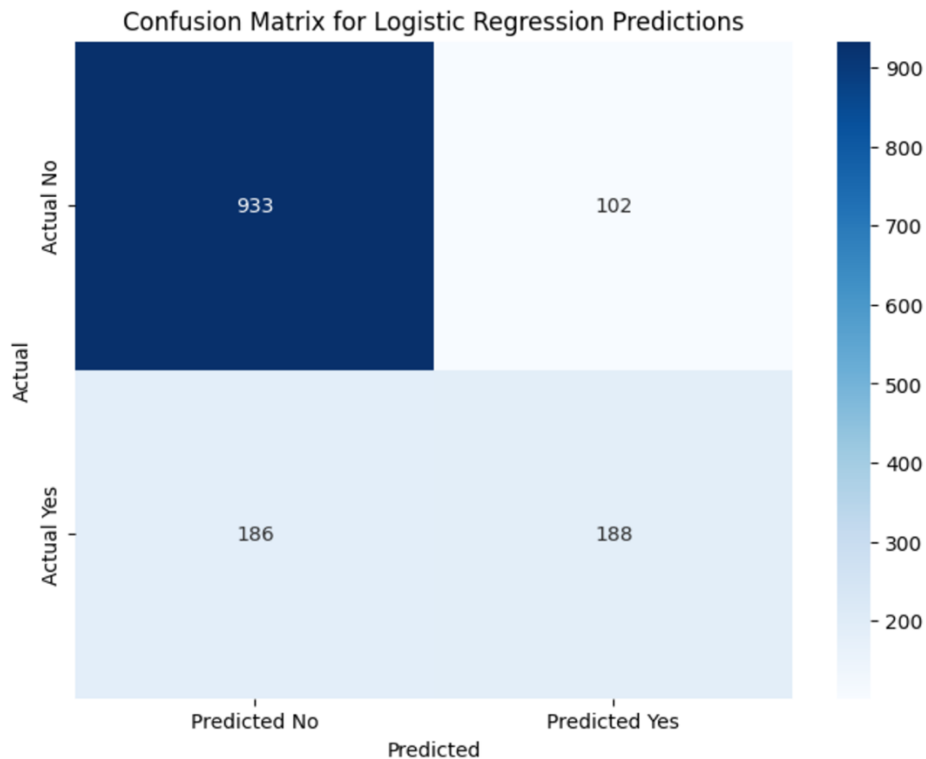
IV. Handle Data Imbalance

## Multinomial Naive Bayes

I. **Feature Preparation**
   - Discretization of **boolean-encoded** or **count-based** features

II. Train Model

III. Hyperparameter Tuning

IV. Handle Data Imbalance

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# Probabilistic Models - Logistic Regression

I. Initialize and Train Model

II. Cross-Validation

III. Predictions and Final Metrics



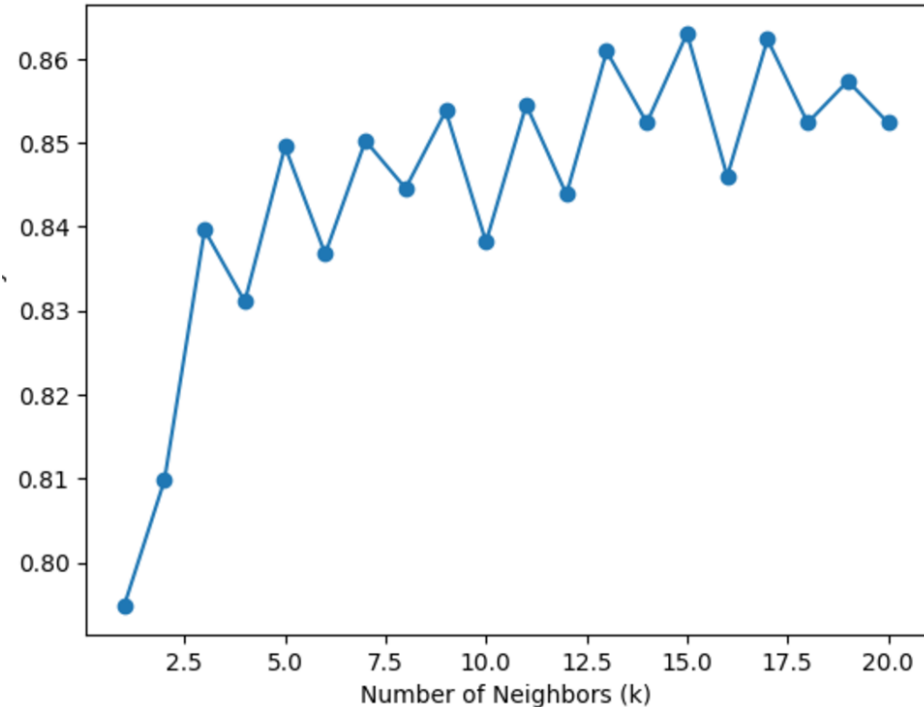Confusion Matrix for Logistic Regression Predictions

# Distance-Based Models - KNN

I. **Standardization**

II. Baseline KNN Model

III. **Find the Optimal Number of Neighbors (k)**

IV. Optimize KNN Model Training

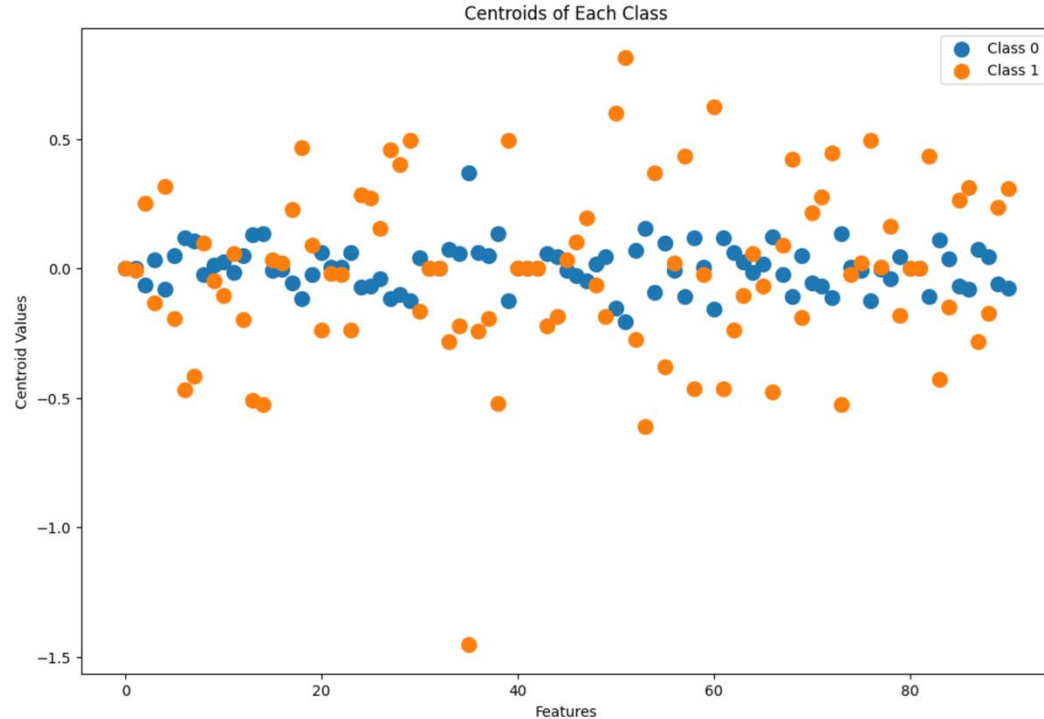V. Final Performance Metrics



KNN Accuracy vs. Number of Neighbors

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

# Distance-Based Models - Nearest Centroid

I. **Standardization**

II. Train Model

III. Cross-Validation

IV. Performance Metrics



Centroids of Each Class

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen
02.12.2024

# Tree-Based Models
## – Decision Tree & Random Forest & XG Boost

## Decision Tree

I.   Data Preparation

II.  **Train Model**
     - Create **entire** Tree

III. Cross-Validation

IV.  Performance Metrics

## Random Forest

I.   Data Preparation

II.  **Train Model**
     - Train a **ensemble** of 100 smaller trees

III. Cross-Validation

IV.  Performance Metrics

## XG Boost

I.   Data Preparation

II.  **Train Model**
     - Build trees **sequentially**

III. Cross-Validation

IV.  Performance Metrics

# Optimization and Kernel-Based Models - SVM & MLP

## SVM

I. **Train Model**
   - Train using the fit() method on the **standardized data**

II. Cross-Validation

III. Performance Metrics

## MLP

I. **Train Model**
   - Combine preprocessing and model training in a **Pipeline**, executed via fit()
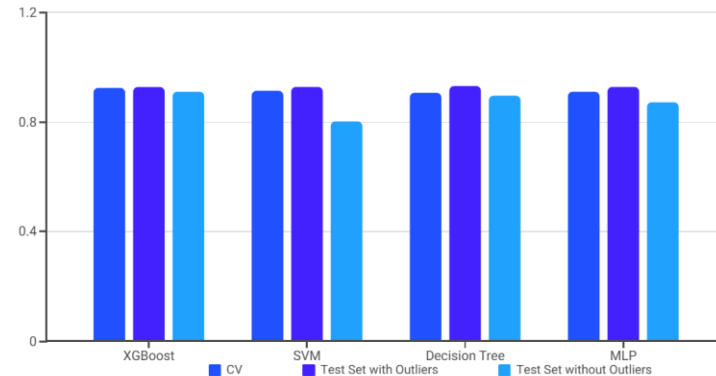
II. Cross-Validation

III. Performance Metrics

# 5. Evaluation

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

02.12.2024

# Evaluation Metrics

### Measurement

- Baseline Model
    - Baseline included predicting the majority class and using rule-based heuristics.
    - Purpose: To understand how much value advanced models add.

- Evaluation
    - Cross Validation Performance
    - Test Set Performance
    - Outlier Remove

# Evaluation Metrics

### Final Results and Reconmendation
- Evaluation based on F1-Score
- XG Boost

| Model | Accuracy | Precision | Recall | F1 | ROC AUC |
|---|---|---|---|---|---|
| XGBoost Prediction_Outlier | 0.954578 | 0.969697 | 0.855615 | 0.909091 | 0.922977 |
| Decision Tree Prediction_Outlier | 0.948900 | 0.966049 | 0.836898 | 0.896848 | 0.913135 |
| MLP Prediction_Outlier | 0.935415 | 0.930091 | 0.818182 | 0.870555 | 0.969774 |
| Random Forest Prediction_Outlier | 0.931157 | 1.000000 | 0.740642 | 0.850998 | 0.870321 |
| SVM Prediction_Outlier | 0.909865 | 0.969582 | 0.681818 | 0.800628 | 0.837044 |
| MNB Prediction_Outlier | 0.789922 | 0.837496 | 0.789922 | 0.800412 | 0.809192 |
| Nearest Centroid Prediction_Outlier | 0.777857 | 0.844765 | 0.777857 | 0.790185 | 0.815494 |
| GNB Prediction_Outlier | 0.845280 | 0.646067 | 0.922460 | 0.759912 | 0.869926 |
| KNN Prediction_Outlier | 0.863023 | 0.829091 | 0.609626 | 0.702619 | 0.782108 |
| Baseline Random Prediction_Outlier | 0.653655 | 0.624612 | 0.653655 | 0.636738 | 0.516647 |
| Baseline Prediction_Outlier | 0.734564 | 0.539584 | 0.734564 | 0.622155 | 0.500000 |
| Baseline Rule Based Prediction_Outlier | 0.598297 | 0.792817 | 0.598297 | 0.614801 | 0.698394 |
| Logistic Regression Prediction_Outlier | 0.795600 | 0.648276 | 0.502674 | 0.566265 | 0.702062 |

# Feature Importance

## Customer Satisfaction is a Primary Driver

The Satisfaction Score stands out as the most significant predictor of customer churn in both SVM and XGBoost models. This aligns with intuition, as customers with low satisfaction are more likely to leave.

## Contract Type Reflects Commitment

Contract types, such as month-to-month and two-year, are highly influential. Month-to-month contracts are associated with higher churn rates while longer-term contracts indicate greater customer commitment.

## Engagement is Crucial

Customer engagement is a strong indicator of loyalty. Features like Number of Referrals and Online Security are highly influential. Customers who are actively engaged with the company are less likely to churn.

## Feature Contributions Vary

The feature importance analysis highlights that different models can assign varying levels of importance to specific features. This emphasizes the need for a holistic understanding of feature contributions across models.

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen

# Preferred Model

## Dashboard

- Predictive performance.

- Interpretability of feature importance.

- Alignment with the dashboard's functionality requirements.

## Preferred Model-XGBoost

- Identify critical features

- Tailor retention strategies

- Enhance customer interactions with precise, data-driven insights.

# Future Work

### Time Series Analysis

Improve churn prediction by incorporating temporal trends in customer behavior. This includes features like changes in usage patterns or payment histories.

### Real-Time Prediction

Integrate real-time data streams into the prediction framework. This would allow for more timely interventions based on changes in customer behavior.

### Dashboard Expansion

Enhance the customer retention dashboard with interactive features. This allows customer service representatives to simulate different retention strategies.

# The Power of Insight: Personalized Recommendations



## Recommendations to Retain Customer:

### Customer Satisfaction
- Low satisfaction. Offer personalized support.

### Security and Support
- Low satisfaction. Offer personalized support.
- No online security. Highlight benefits.

### Contracts and Payments
- Month-to-month contract. Promote long-term benefits.
- Not on two-year contract. Discuss perks.
- No credit card payment. Recommend for convenience.

### Streaming and Services
- Young customer. Highlight appealing services.

### Referrals and Offers
- Low satisfaction. Offer personalized support.

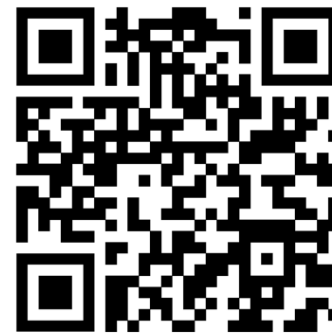# Great, to keep you as a customer!

Thank you for your attention!

github.com/eelisee/telco-customer-churn

Team 9: Zhiqi Yang, Elise Wolf, Xi Liu, Shiqi Zhou, Yenan Chen