



# THE HARTFORD PROJECT

**GROUP 3**

SALISA ALMEIDA

PETER HOFSTEDT

JATIN BOMRASIPETA

# PRESENTATION AGENDA

## 1

- Project Overview
- Business Problems
- Key Insights

## 2

- Data overview
- Challenges on Modeling

## 3

- Two Step Modeling Approach
- Step 1
- Step 2

## 4

- Customer Segmentation and Retention Strategies
- Next Step

# PROJECT OVERVIEW



## The primary objective

- Analyze factors influencing customer retention
- Develop predictive models that can help identify at-risk customers
- Set strategies ideas to improve retention rates.

# BUSINESS PROBLEM



Can the likelihood of a customer canceling their auto insurance be estimated?



Additionally, can we determine the specific type of cancellation they will choose, such as midterm or flat cancellation?

# KEY INSIGHTS

## MODEL

**Performance:** LightGBM showed the highest accuracy on Test set.  
Precision = 0.84 for renewals

0.82 for cancellations  
AUC = 0.705

**Feature Engineering:** Created 5 new features that capture the data better and reduce dimensionality for improved model performance.

## AS BUSINESS

**Customer Segmentation:** Identified four distinct customer segments with varying loyalty levels and retention needs.

**Targeted Retention Strategies:** Developed tailored strategies for each segment to improve retention rates and customer satisfaction.

# DATA OVERVIEW AND MODEL CHALLENGES

1

## Data Cleaning and Enhancement

- Before Cleaning: 264,963 rows by 54 Columns
- After Cleaning: 264,963 rows by 57 Columns
- Removed 2 variables, created 5 variables

2

## Preparing for Modeling

- Label Encoding
- Checking correlation between variables
- Checking correlation with STATUS

3

## Unbalanced Classes:

- STATUS:
  - 0 = Not Cancelling - 84.6%
  - 1 = Mid-Term Cancellation - 8.2%
  - 2 = End-Term Cancellation - 7.3%

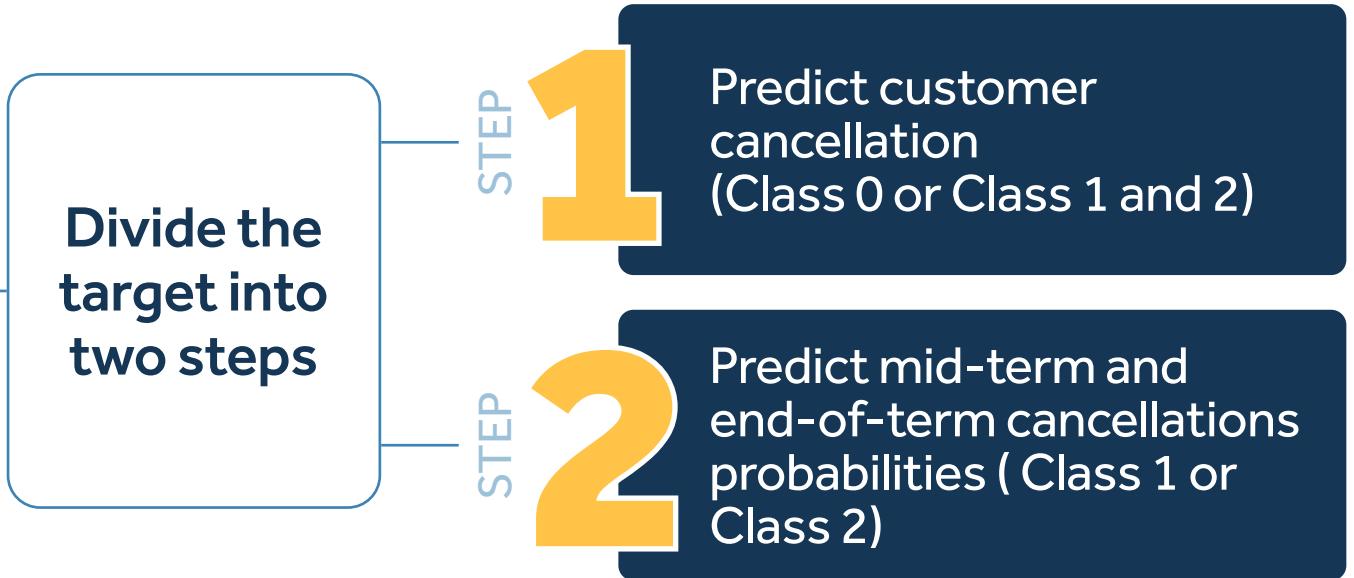
4

## Balancing Techniques

- SMOTE
- Downsampling



# MODELING APPROACH



**Metrics:** We focused on precision because it helps accurately identify real cancellations and minimize false positives in predicting customer cancellations. ROC-AUC was used too as the data imbalance was severe

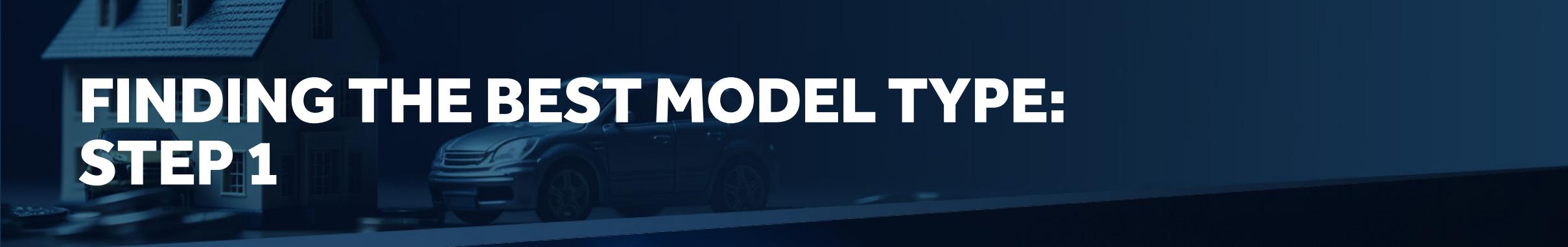
# CHOOSING OUR FEATURE SELECTION METHOD FOR STEP 1



- Recursive Feature Elimination
- Principal Component Analysis
- Feature Importance ranking
- No classifier Feature selection
  - Correlation with STATUS thresholds and correlations with other variables

The best method for feature selection that we found was **Feature Importance Ranking**

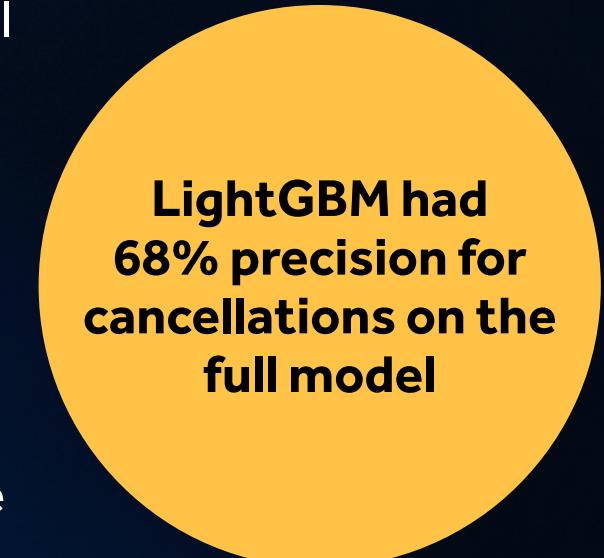
# FINDING THE BEST MODEL TYPE: STEP 1



We tested different models on the full dataset to choose the best model

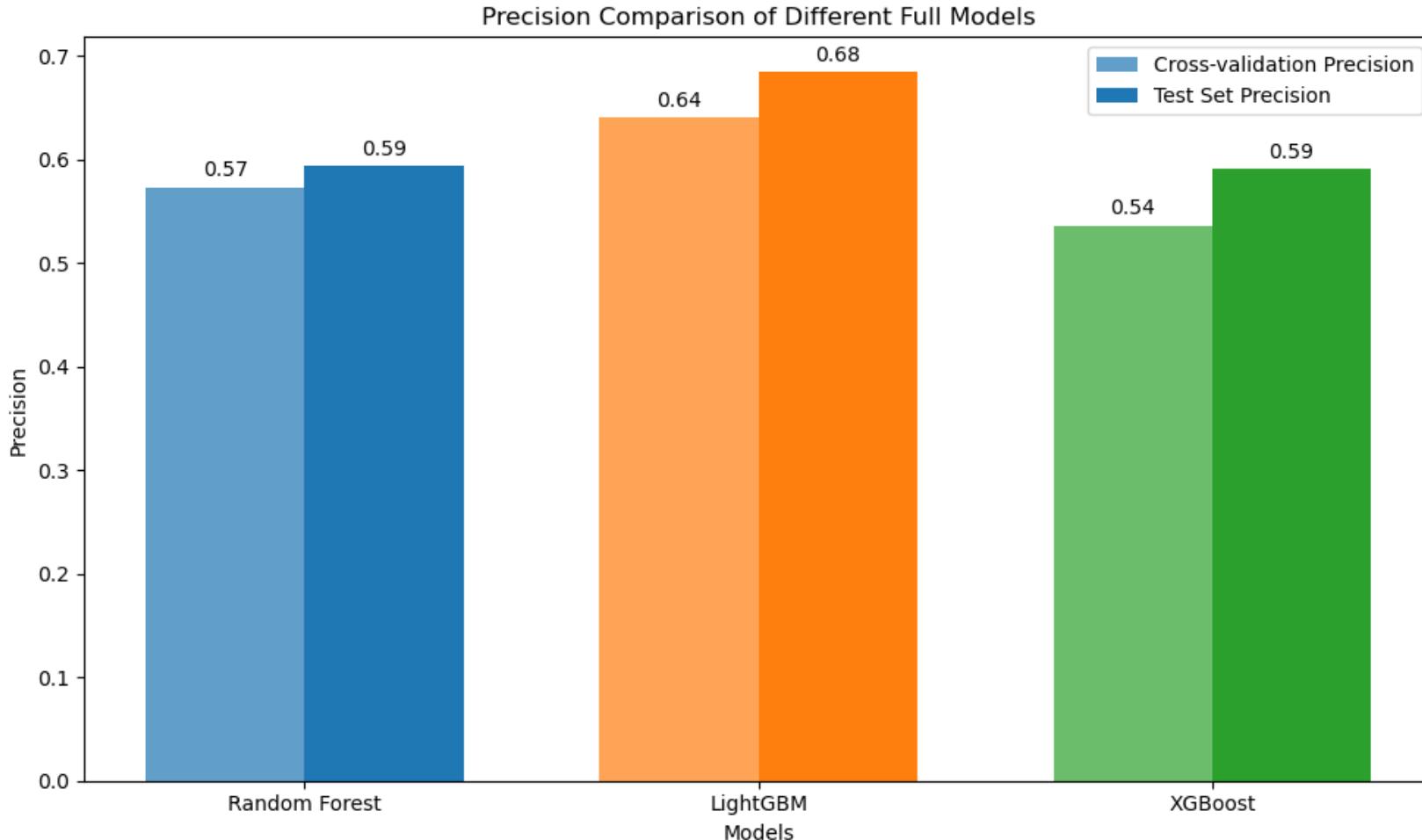
- Random Forest
- XGBoost
- LightGBM

We found that **LightGBM** had overall higher accuracy and better predictive power than the other models



**LightGBM had  
68% precision for  
cancellations on the  
full model**

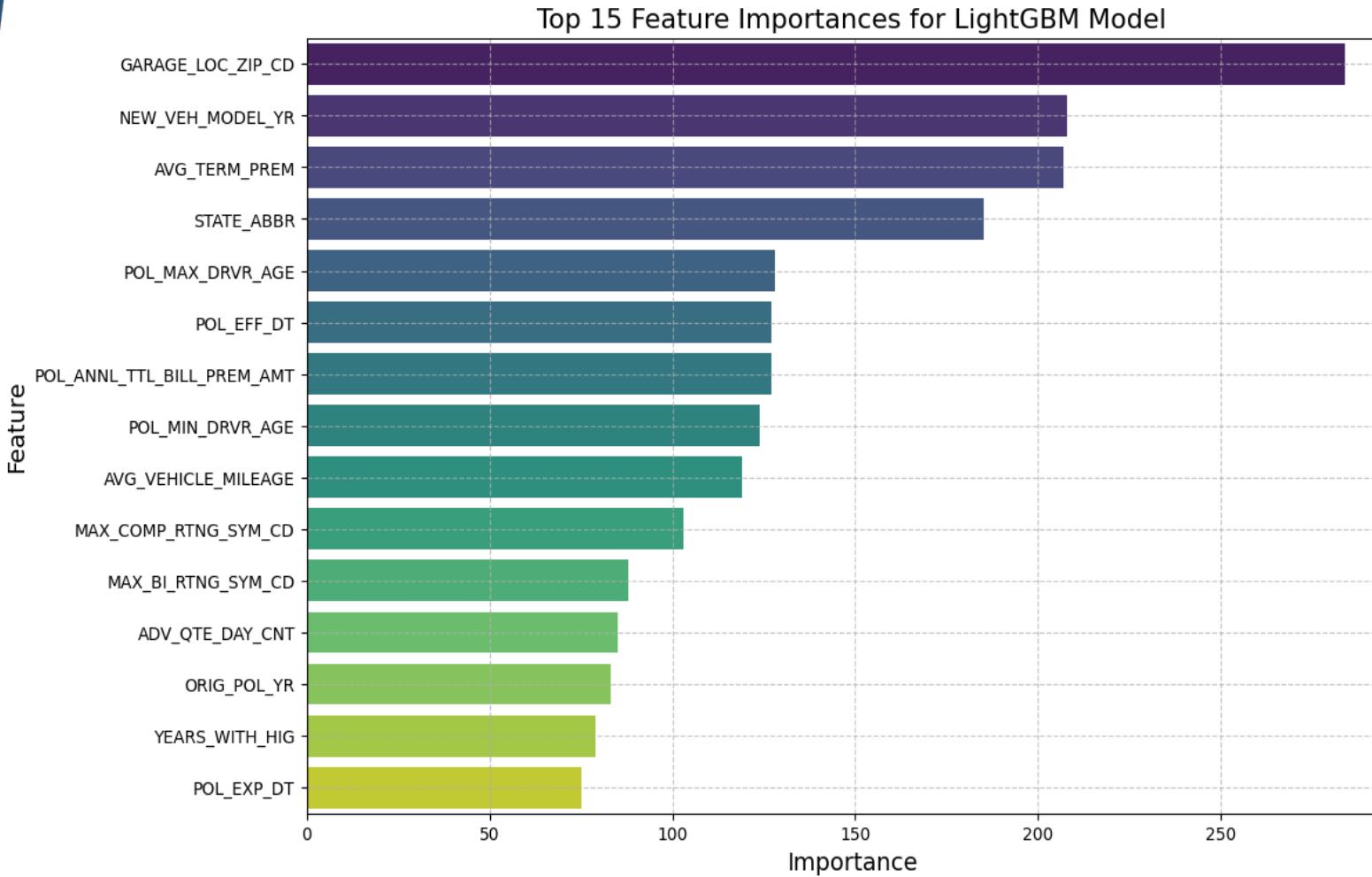
# COMPARISON OF DIFFERENT STEP 1 MODELS ON THE FULL DATASET



LightGBM had a better performance using all variables.

CV = 0.64  
Test = 0.68

# FEATURE IMPORTANCES

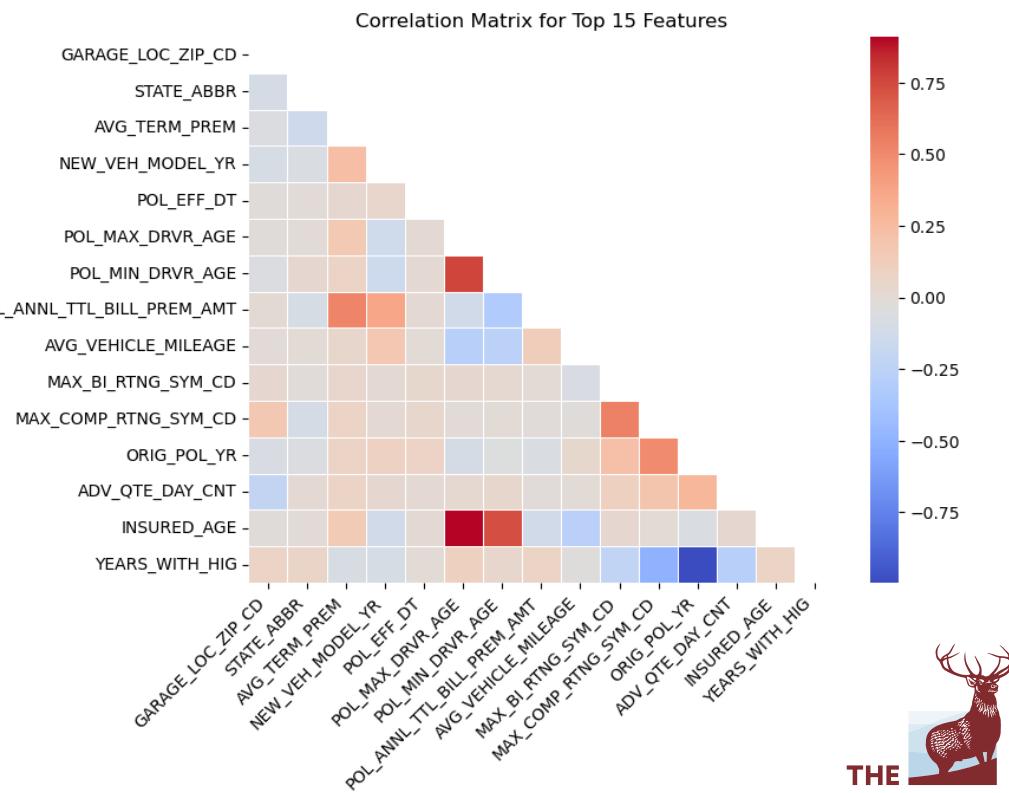
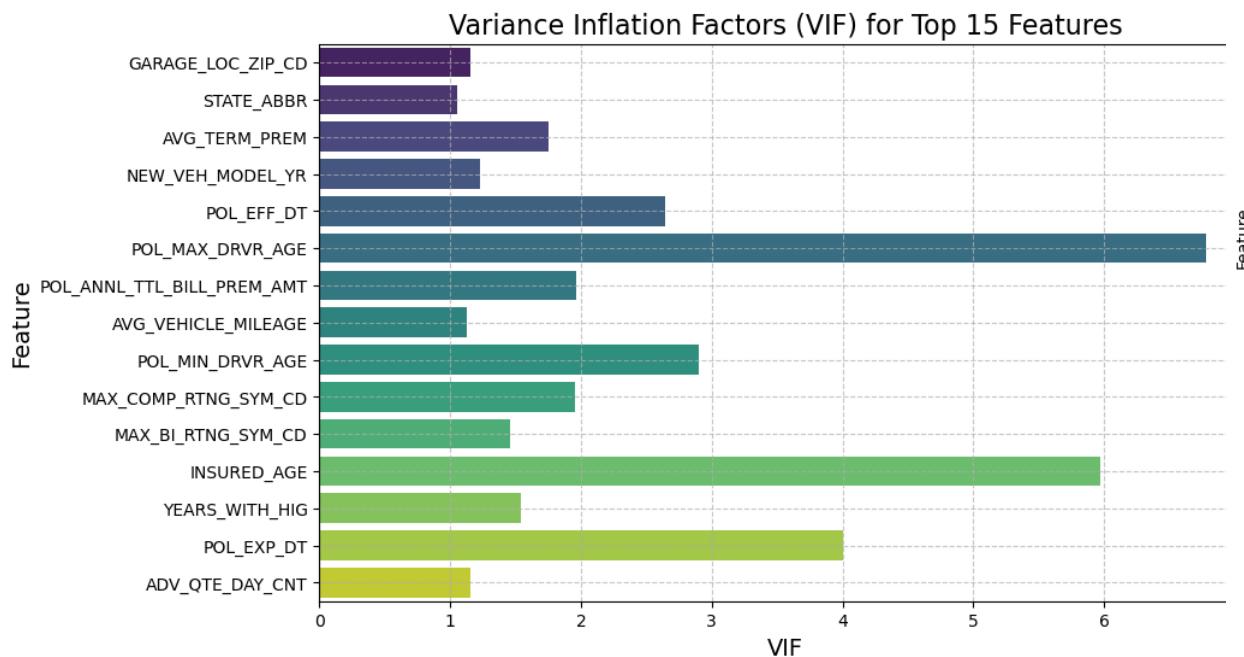


LightGBM  
feature  
selection  
technique  
was used for  
selecting the  
best 15

# CORRELATION AND VIF

## 15 Best Variables Cleaning

- ORIG\_POL\_YR was removed for correlation.
- POL\_EFF\_DT and POL\_EXP\_DT Removed for VIF and high correlation.



# DATA ENHANCEMENT

- Create new variables that show important business statistics.
- Increased robustness
- Reduce the complexity of models and improve efficiency
- Ensure model accuracy

# FEATURE ENGINEERING FOR ENHANCED MODEL PERFORMANCE

## NEW FEATURES

<b>PREMIUM_CHANGE:</b> Absolute change in renewal premium over two terms.	<b>DRV_R_AGE_DIFF:</b> Difference between the oldest and youngest drivers on the policy.	<b>BILL_PREM_DIFF:</b> Difference between the total annual bill premium amount and the average term premium.	<b>AVG_COMP_BI_RTNG_SYM:</b> Average of comprehensive and bodily injury rating symbols.	<b>VEHICLE_AGE:</b> Age of the vehicle calculated from the model year.

Variables used to create new features were then removed to avoid multicollinearity

# STEP 1 FINAL MODEL FEATURE SELECTION

## Original Features

BILL\_PRE\_DIFF  
VEHICLE\_AGE  
YEARS\_WITH\_HIG  
STATE\_ABBR  
GARAGE\_LOC\_ZIP\_CD

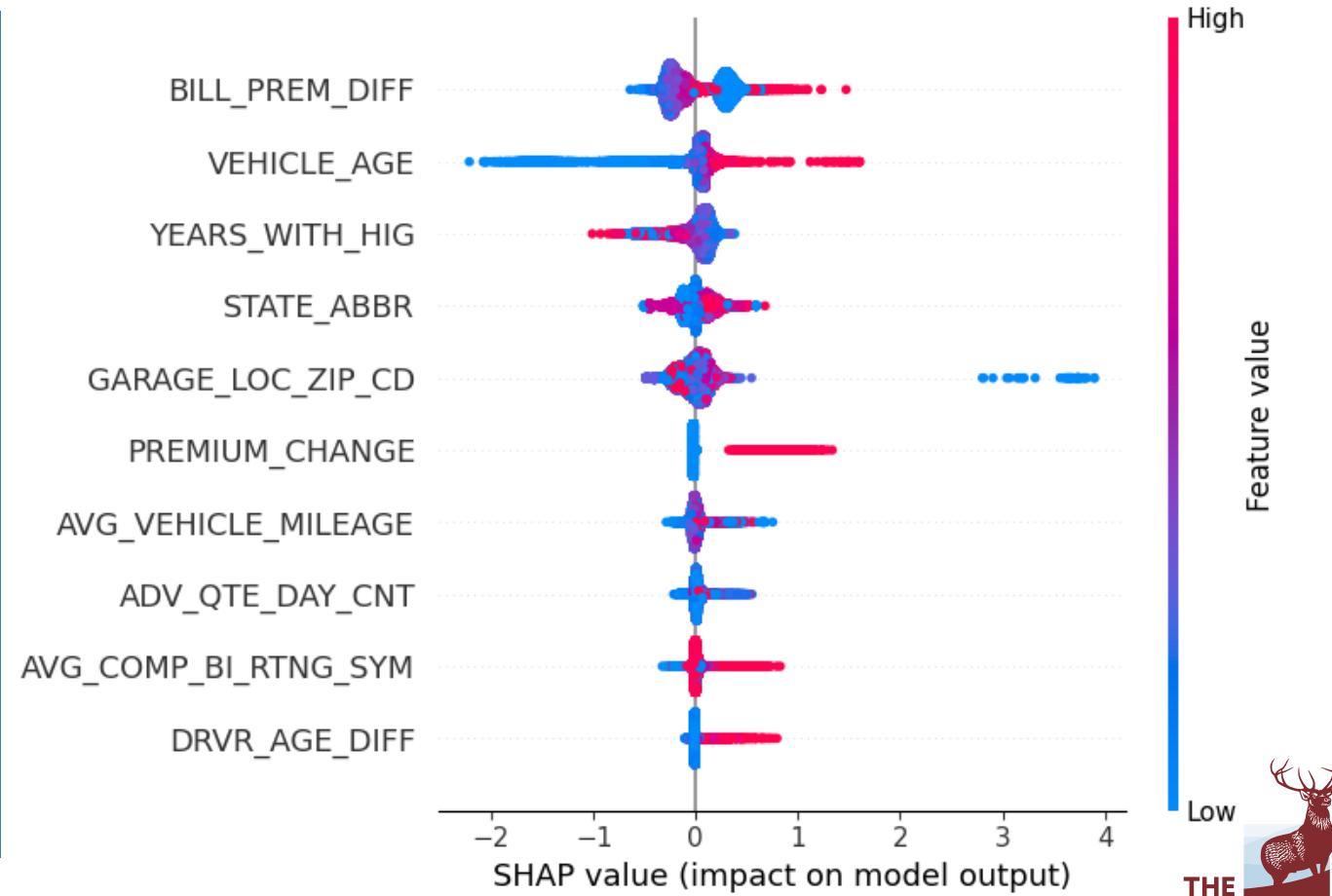
## New Features

PREMIUM\_CHANGE  
AVG\_VEHICLE\_MILEAGE  
ADV\_QTE\_DAY\_CNT  
AVG\_COMP\_BL\_RTNG\_SYM  
DRV\_R\_AGE\_DIFF

# CHECKING SHAP VALUES FOR THE FIRST MODEL RESULTS

Positive SHAP values indicate an increase in the predictive power of cancellation and vice-versa.

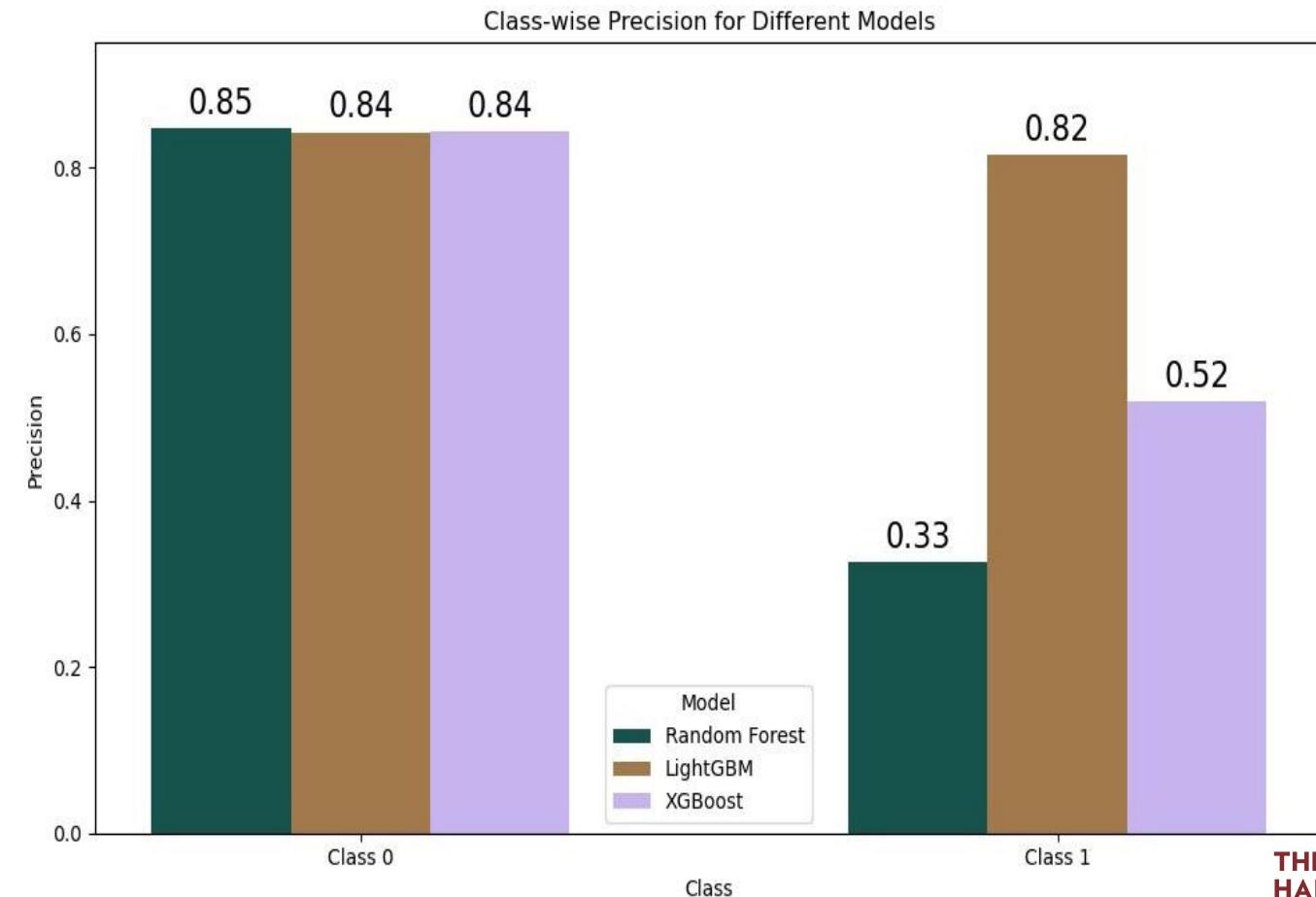
SHAP Values best features for predicting cancellations:  
**BILL\_PREM\_DIFF**,  
**VEHICLE\_AGE** and  
**YEARS\_WITH\_HIG**



# STEP 1 FINAL MODEL PERFORMANCE

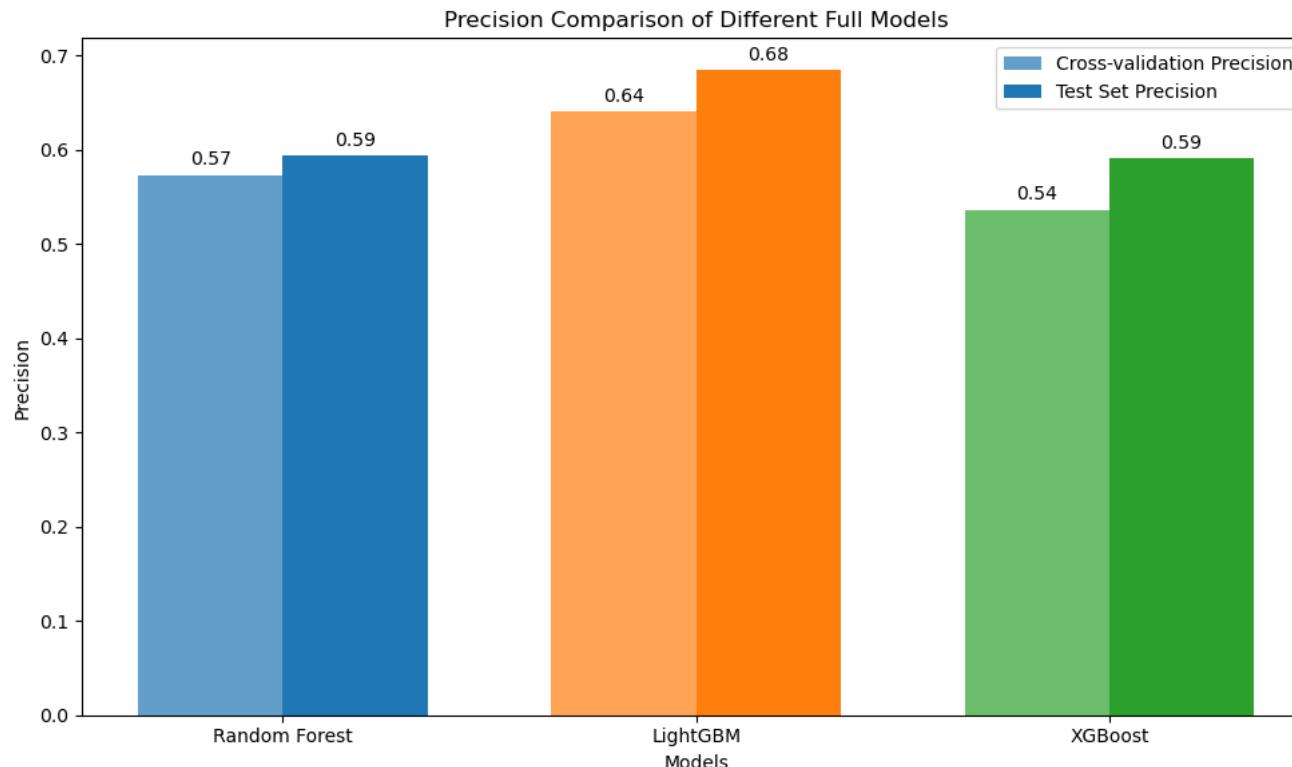
Step 1 Final model  
Performance  
comparison of  
cancellations  
between models.

LightGBM results in:  
Precision class  
 $0 = 84, 1 = 82$   
ROC-AUC = 0.70

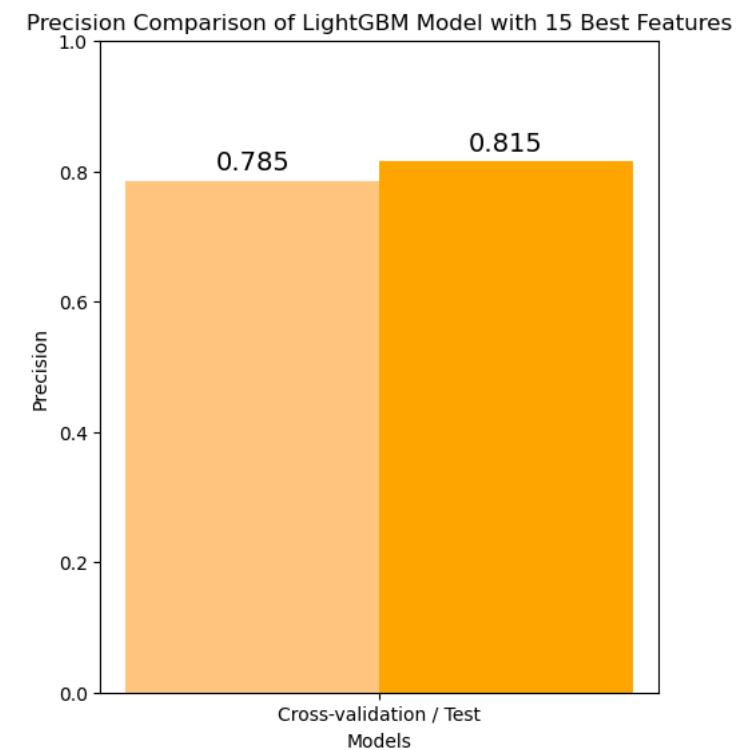


# STEP 1 MODEL RESULTS ON THE TEST SET

## Full dataset prediction



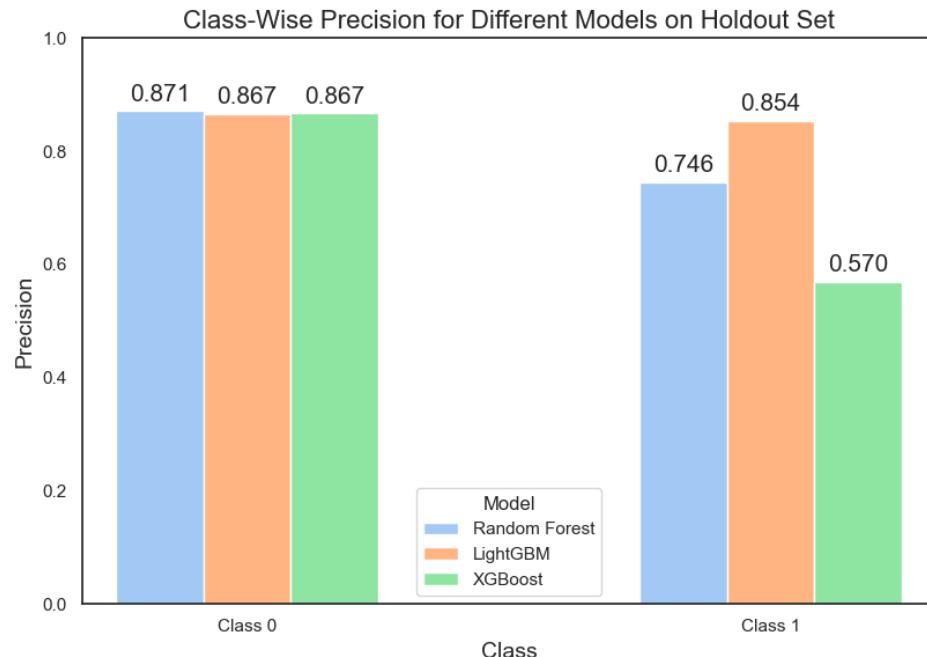
## Step 1 Final Model Prediction



LightGBM had a better performance using all variables and reached 0.82 on Test set using the 15 best variables for prediction.

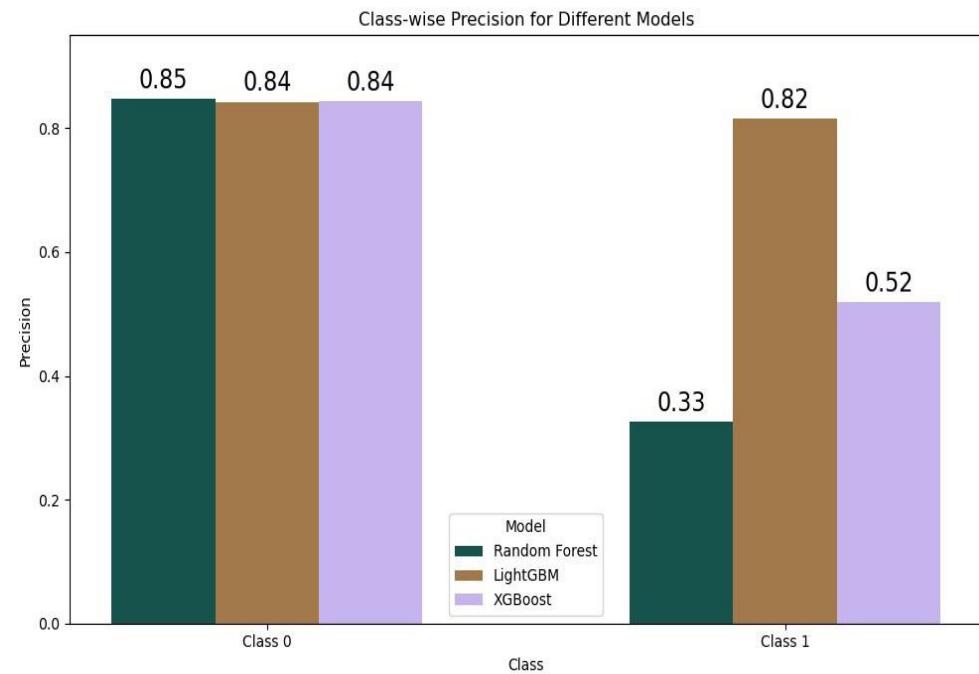
# STEP 1 MODEL - HOLDOUT SET VS TEST SET

## HOLDOUT RESULTS



The consistent performance across the test set and holdout set indicates that the model is not overfitting. ROC-AUC = 0.65

## TEST RESULTS



The model showed a slightly better performance on the holdout set.  
ROC-AUC = 0.70

# STEP 1: OTHER APPROACHES

Techniques we tried to implement to improve our prediction results

## Weighting

Increases other metrics at the cost of precision

## SMOTE

No changes on classes precision

Lasso and Ridge regularization  
Did not change precision

## PCA

Class 1 was overestimated

## Learning Rate Scheduler

# STEP 2 MODEL PREDICTION

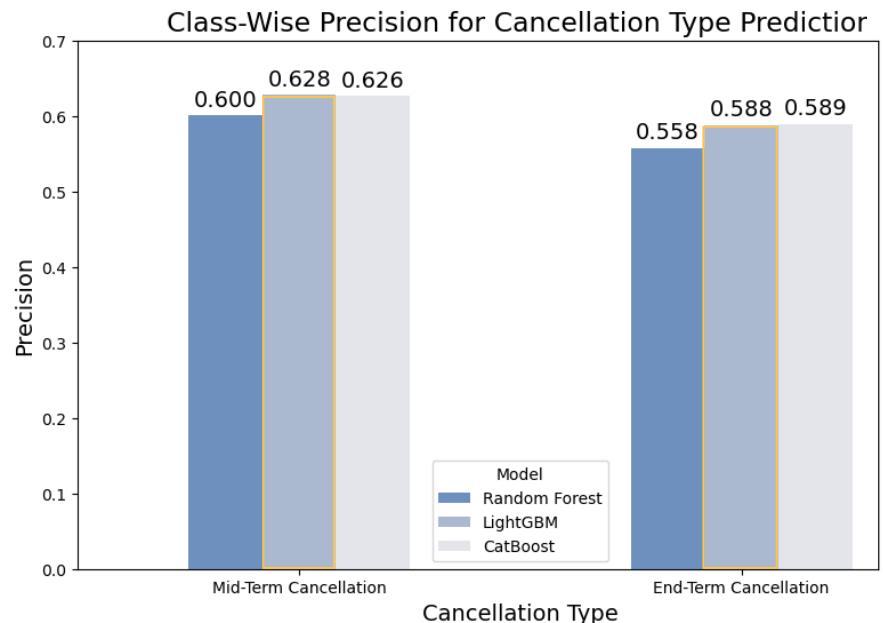
**Target Variable Creation:**  
Filtered out non-cancellations and used STATUS column for cancellation type.

**Feature Selection:**  
Used LightGBM, identified ZIP Code and Policy Effective Date as top features. Various feature selection techniques were applied, but none significantly improved model performance.

**Challenges Faced:**  
Insufficient data for minority classes, impacting performance despite various techniques.

STEP 2 Cancellation type model: Best results were achieved by 15 top feature importance selection using LightGBM classifier

## STEP 2 Model Results



Precision didn't meet levels we found acceptable: Class 1: 63% Class 2: 59%



# CUSTOMER SEGMENTATION OVERVIEW

We focused on **3 KEY FEATURES** to divide the customers into **4 SEGMENTS**  
**'BILL\_PREM\_DIFF', 'VEHICLE\_AGE' AND 'YEARS\_WITH\_HIG'**

METHOD USED:  
KMEANS  
CLUSTERING

	Premium Difference	Vehicle Age	Avg Years with HIG
SEGMENT 0	Moderate	Older Vehicles	7+
SEGMENT 1	Very Low	Oldest Vehicles	6+
SEGMENT 2	Highest	Newest Vehicles	9+
SEGMENT 3	Relatively High	Relatively New	8+

# CUSTOMER SEGMENTATION VISUALS

	SIZE	CUSTOMER LOYALTY	STRATEGIES
SEGMENT 0	100,000+	HIGH (LONG TERM CUSTOMERS)	LOYALTY PROGRAMS & CUSTOMER FEEDBACK
SEGMENT 1	100,000+	MODERATE (POTENTIAL FOR GROWTH)	ENGAGEMENT CAMPAIGNS & INCENTIVES ON RENEWALS
SEGMENT 2	ABOUT 5,800	VERY HIGH	PREMIUM SERVICES & PERSONALIZED OFFERS
SEGMENT 3	AROUND 40,000	HIGH BUT NEEDS ENGAGEMENT	REGULAR CHECK INS & TARGETED DISCOUNTS

# NEXT STEPS

**1**

**Gather More Data:**  
Collect additional information on customers who cancel mid-term to better understand their characteristics and reasons for cancellation.

**2**

**Integrate External Data Sources:** Incorporate external data sources, such as economic indicators and competitor actions, to enhance the model's predictive power and provide a more comprehensive view of the factors influencing customer cancellations.

**3**

**Continue the Development of Step 2 Model:**  
Find and implement new techniques that could help increase the predictive power of the step 2 model.

The background of the slide is a dark blue color. In the center, there is a small, semi-transparent image of a two-story house with a chimney and a gabled roof. To the left of the house, a dark-colored SUV is parked. To the right, a smaller car is partially visible. Scattered around the base of the house are several silver coins of different sizes.

Thank you

