# Table of Contents

# 1. Executive Summary

## 1.1 Study Proposal

This study aims to apply machine learning techniques, and business research to understand and address low engagement rates in preventive care visits among Humana's Medicare Advantage Local Preferred Provider Organization (LPPO) members. Despite the flexibility offered in specialist access, LPPO plans experience significantly higher disengagement compared to HMO plans.

We will provide Humana with a detailed understanding of the factors driving disengagement in preventive care among LPPO members. The combination of predictive modeling and qualitative research will lead to practical, data-driven strategies aimed at increasing member engagement. If implemented effectively, these strategies will improve Humana's Stars ratings, ensure more accurate risk documentation, and enhance overall financial performance while benefiting member health outcomes.

## 1.2 Modeling and Analysis

Our analysis aimed to predict which LPPO members are at the highest risk of disengaging from preventive care visits, focusing on the disengagement outcome as the target variable. We utilized a comprehensive training dataset that included demographic information, healthcare utilization patterns, and claims data. For this task, we compared different machine learning models and LightGBM was selected due to its superior performance in handling large datasets and imbalanced classes, making it an ideal choice for our predictive modeling.

Furthermore, we refined the model using a combination of Bayesian optimization and Randomized Search for hyperparameter tuning. Bayesian optimization allowed us to efficiently explore the hyperparameter space by building a probabilistic model and focusing on the most promising areas, thus reducing the number of trials needed to find optimal parameters. This method is particularly effective for models with a complex search space, as it balances exploration and exploitation of hyperparameters.

In parallel, Randomized Search was employed to quickly test a wide range of hyperparameter combinations, providing flexibility and ensuring a diverse set of configurations were evaluated. By combining both approaches, we achieved more robust and reliable tuning, enhancing the model's performance by improving accuracy, reducing overfitting, and optimizing resource usage during training. This dual strategy allowed us to fully exploit the potential of our model and increase its generalizability on unseen data.

## 1.3 Results and Recommendations

Based on our analysis, we identified four key strategies for Humana to adopt in order to improve LPPO member engagement in preventive care:

1. Reducing financial strain through deductible reductions and partial co-insurance waivers.
2. Enhancing access to care by expanding home-based care services and strengthening virtual care options.
3. Streamlining administrative processes with improved digital tools and support for claims submission.
4. Providing personalized member support through AI-driven case managers and structured support groups.

These data-driven recommendations will help Humana make informed decisions to increase member engagement, improve health outcomes, and enhance overall satisfaction, ultimately ensuring better performance in the Medicare Advantage market.

## 2. Case Background

This case revolves around Humana's Medicare Advantage (MA) program, specifically focusing on the engagement challenges associated with Local Preferred Provider Organization (LPPO) plans. While LPPOs offer members greater freedom in selecting care providers compared to Health Maintenance Organizations (HMOs), they face a higher proportion of unengaged members. These unengaged members do not utilize preventive primary care visits with their Primary Care Physicians (PCPs), which are vital for managing chronic diseases, conducting preventive health screenings, and assessing health risks, especially among aging populations.

Medicare Advantage is structured to provide comprehensive coverage by bundling hospital, outpatient, and prescription drug benefits. LPPO plans provide the flexibility of direct specialist access without a PCP referral, but this flexibility seems to come at the cost of reduced engagement in preventive care. This disengagement leads to fewer touchpoints with healthcare providers, potentially lowering Stars performance (Medicare's quality rating system) and missing opportunities to properly document health risks.

The Centers for Medicare & Medicaid Services (CMS) use Medicare Risk Adjustment (MRA) to ensure that MA organizations are compensated based on the health status of their members. Inaccurate or incomplete documentation, particularly among unengaged members, threatens the financial performance of LPPO plans. The decrease in preventive visits affects Humana's ability to maintain or improve its Stars ratings and risks reducing the bonuses that would otherwise be invested in member benefits.

## 2.1 Business problem

This analysis aims to address a significant business challenge faced by Humana's Medicare Advantage LPPO plans: the low engagement rates for preventive care visits with Primary Care Physicians (PCPs). Despite offering enhanced flexibility in accessing specialists without referrals, LPPO plans exhibit a markedly higher proportion of unengaged members compared to Health Maintenance Organization (HMO) plans.

Unengaged members fail to utilize preventive care visits, which are essential for the management of chronic conditions and for ensuring accurate health documentation. Based on our findings, we identified two primary challenges associated with member disengagement for Humana:

1. **Decline in Stars Performance**
   Reduced engagement with PCPs leads to fewer preventive health screenings and subsequently lower Stars quality ratings. This decline adversely impacts Humana's capacity to secure performance bonuses that fund additional member benefits.
2. **Incomplete Risk Documentation**
   A lack of preventive visits hinders the precise documentation of members' health risks, negatively affecting the Medicare Risk Adjustment (MRA) payments that are critical for providing care to high-risk patients.

Consequently, the objectives of this analysis are two fold:

1. To identify key characteristics and behaviors that predict which LPPO members are at the highest risk of disengagement from preventive care.
2. To develop actionable insights and strategies through data modeling and business research to enhance member engagement and improve overall performance.

## 2.2 Key Performance Indicators

### A. AUC-ROC Curve

The AUC-ROC (Area Under the Receiver Operating Characteristic) Curve is a crucial metric for evaluating the performance of our classification model, which predicts member engagement with preventive care in LPPO plans. The ROC curve plots the True Positive Rate (sensitivity) against the False Positive Rate (1-specificity) across various thresholds, providing a measure of the model's discriminatory ability.

In this context, a higher AUC indicates better model performance in distinguishing between engaged and unengaged members. A model with an AUC closer to 1 implies it is more accurate in identifying members who are likely to remain unengaged and those who are engaged. Given the impact of preventive care on health outcomes and Humana's Stars performance, we aim for a high

AUC to ensure that the model can effectively predict unengaged members and guide intervention strategies.

## B.  AUC for Precision-Recall Curve

Since unengaged members may represent a smaller proportion of the overall population, the Precision-Recall (PR) curve is a critical tool for evaluating model performance on this imbalanced data. The PR curve focuses on the minority class (unengaged members) by plotting Precision (the fraction of predicted unengaged members that are actually unengaged) against Recall (the fraction of actual unengaged members correctly identified).

The area under the PR curve provides a summary of the model's effectiveness across different levels of Precision and Recall. Our primary goal is to optimize the model for accurately predicting unengaged members, as the cost of failing to identify these members (false negatives) could lead to missed opportunities for preventive care outreach. Ideally, we aim for an AUC value close to 1 for the PR curve.

## C. Confusion Matrix:

The confusion matrix gives a detailed breakdown of the model's predictions, categorized into four outcomes: True Positives (correctly predicted unengaged members), True Negatives (correctly predicted engaged members), False Positives (engaged members incorrectly predicted as unengaged), and False Negatives (unengaged members incorrectly predicted as engaged).

For Humana's LPPO plans, minimizing False Negatives is crucial since they represent missed opportunities to engage members in preventive care. At the same time, we aim to maximize True Positives to correctly identify unengaged members who would benefit from intervention. The confusion matrix provides valuable insights into where the model might need further refinement to achieve a balance between correctly predicting unengaged and engaged members.
Ideal Scenarios:
• True Positives: High, as they represent unengaged members correctly identified for outreach.
• True Negatives: High, representing engaged members accurately predicted to stay engaged.
• False Positives: Low, to reduce unnecessary outreach efforts to engaged members.
• False Negatives: Minimized, as these represent missed chances to engage unengaged members.

## D.  Precision, Recall, and F1 Score

In this business problem, high Recall (sensitivity) is prioritized to capture as many unengaged members as possible. Precision ensures that the unengaged members identified are actually unengaged, preventing misallocation of resources. The F1 Score balances Precision and Recall, providing a single measure that reflects both, making it an important metric in guiding outreach strategies to improve engagement.

# 3. Data Preparation

## 3.1 Data Overview

There are two datasets provided by Humana this year, the training dataset, the one used to train the model, and the holdout dataset, the one used to make predictions. The target variable is identified as *'preventive_visit_gap_ind'*.

### 3.1.1 Target Variable

The target variable can take values of either 0 or 1, and the below visualization gives the therapy scenarios for each of the two values for this flag indicator. The training dataset includes 1,527,904 Humana members' information. The holdout dataset includes 381,976 Humana members' information.

### 3.1.2 Datasets Descriptions

1. **Demographics:** Contains key characteristics such as age, sex, race, region, and member status, including tenure with the plan and plan type. These features provide a foundational understanding of the member population and their demographic distribution.
2. **Claims Data:** Provides comprehensive information about healthcare service utilization, pharmacy usage, claims history, and costs. This includes data on inpatient admissions, outpatient visits, and prescription drug usage, offering insights into how members interact with healthcare services over time.
3. **Social Determinants of Health:** Includes various social-economic indicators that can influence member engagement and health outcomes. These metrics cover poverty rate, unemployment rate, healthcare access, smoking rates, obesity levels, and other behavioral risk factors that may affect overall health and healthcare utilization.
4. **Web Activity:** Tracks member engagement with Humana's digital platforms, capturing details such as login frequency and the number of days since the last login. This data can be used to understand digital engagement trends and member interaction with online health management tools.
5. **Pharmacy Utilization:** Captures data related to prescription drug usage, including details on the types of drugs prescribed, costs, and the frequency of prescription fills. This section offers insights into member reliance on medication as part of their healthcare management.
6. **Cost & Utilization:** Provides insights into overall healthcare spending and claims utilization, broken down by service type. It includes distinctions between in-network and

out-of-network services, helping to analyze patterns in healthcare resource usage and cost efficiency.

7. **Marketing Control Point:** Tracks interactions between members and healthcare providers through various communication channels. This includes the number of emails, phone calls, and live calls, helping to assess how effectively members are being engaged through direct communication strategies.

8. **Member Condition:** Contains data on members' chronic conditions and health status, including chronic condition classifications and indices such as the diabetes severity index. This data is crucial for understanding the health challenges faced by members and tailoring healthcare interventions accordingly.

9. **Member Details:** Includes detailed demographic attributes, such as veteran status, geographic location, and tenure with the healthcare plan. This data helps provide context for member health behaviors and healthcare needs.

10. **Member Claim History:** Offers information about healthcare interactions, including primary care physician (PCP) visits, specialist consultations, preventive care visits, and telehealth engagements. This history is key to understanding how members access and utilize various forms of healthcare.

## 3.2 Data Exploration

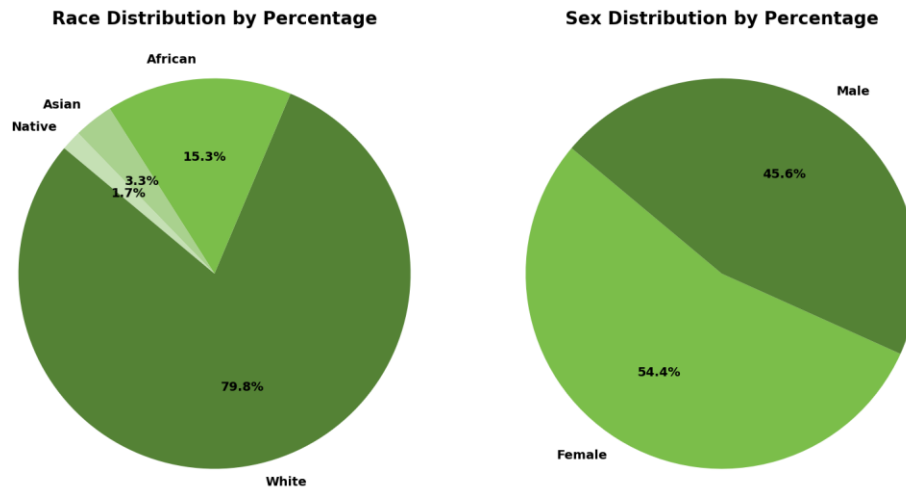### 3.2.1 Exploratory Data Analysis

#### (1) Race Distribution

The chart shows that the majority of the population is White (79.8%), followed by African (15.3%), Asian (3.3%), and Native (1.7%). This indicates that White individuals make up the largest portion of the dataset, and understanding the healthcare behaviors and needs of these different racial groups is important for targeted interventions.

#### (2) Sex Distribution

The chart reveals a nearly even split between females and males, with 54.4% being Female and 45.6% Male. This gender balance provides an opportunity to analyze how men and women engage with healthcare services and whether there are differences in preventive care utilization between the two groups.
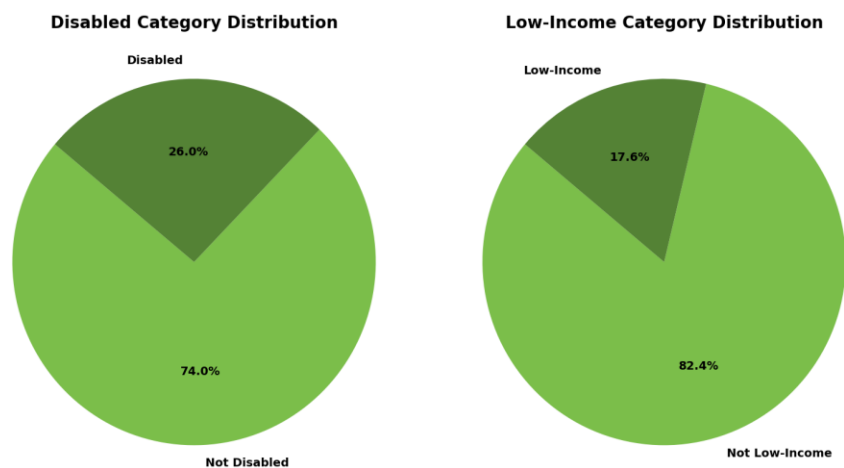
**Figure1. Race and Sex Distribution**

## (3) Disabled Category Distribution

The chart, 26.0% of the individuals are classified as disabled, while 74.0% are not disabled. This highlights the importance of considering the needs of disabled individuals, who may require more specialized healthcare services and face additional challenges in accessing care.

## (4) Low-Income Category Distribution

The chart shows that 17.6% of the population is classified as low-income, while 82.4% are not. This indicates that a significant minority of the population may face financial barriers that could impact their healthcare access and engagement, making it a key consideration for any interventions aimed at increasing preventive care utilization.

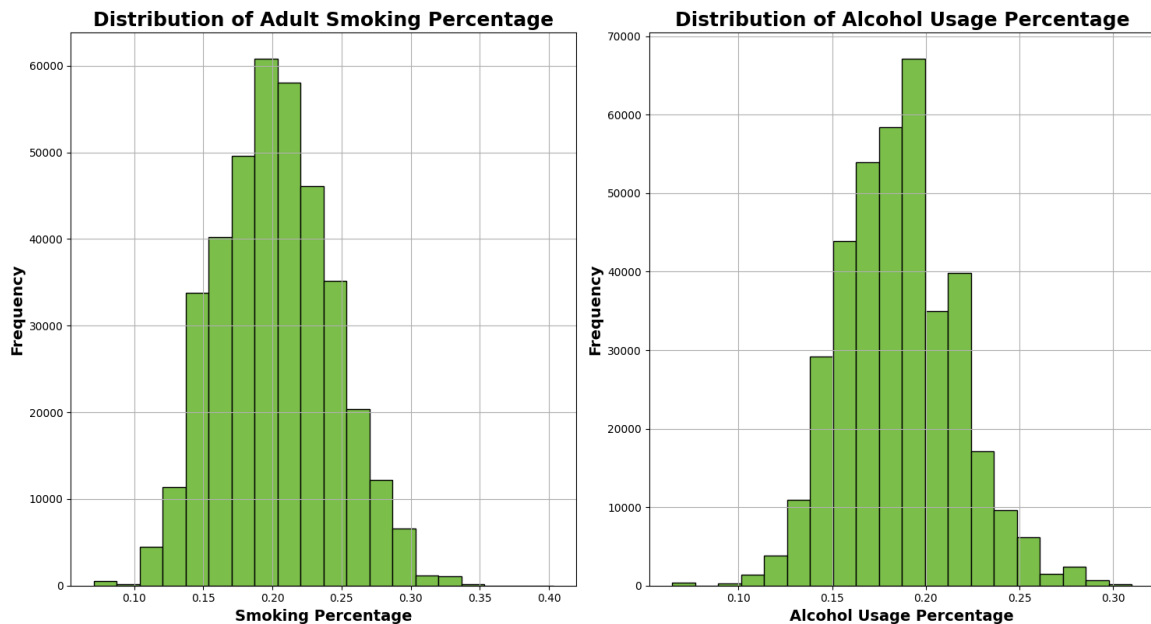**Figure2.  Disabled Category and Low-Income Distribution**

### (5) Adult Smoking Distribution

The distribution follows a fairly symmetric bell-shaped curve, indicative of a normal distribution, with its peak around the 20% mark. This pattern suggests that a significant portion of the population experiences smoking rates near this percentage. The majority of smoking rates range between 10% and 35%, with only a few areas showing values below or above this range. The most frequent smoking rate is approximately 20%, indicating that, on average, one in five adults in the observed regions smokes.

### (6) Alcohol Usage Distribution

The distribution also exhibits a bell-shaped pattern, peaking around 20%, much like smoking rates. The alcohol usage rates predominantly fall between 10% and 30%, with minimal instances of rates below 10% or above 30%. The most common alcohol usage rate also hovers around 20%, suggesting that, on average, one in five adults in these regions regularly consumes alcohol.

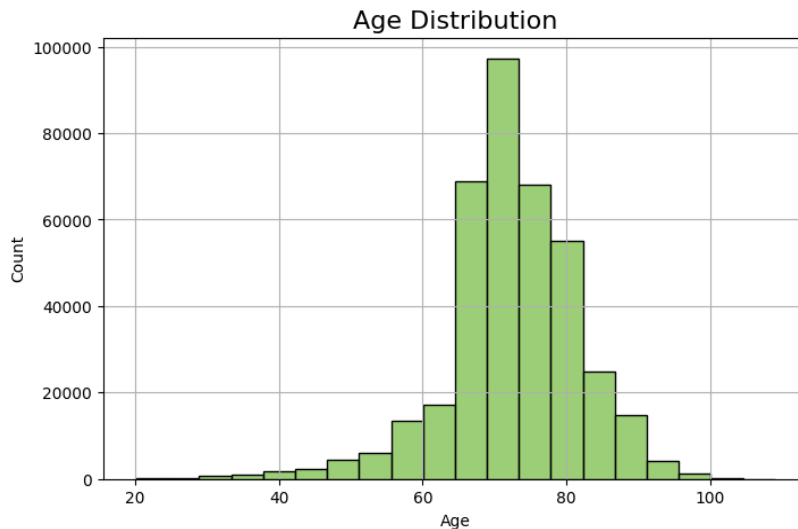Both smoking and alcohol usage distributions demonstrate similar patterns, with their highest frequencies clustering around 20%. This indicates a consistency in these lifestyle risk factors across different regions in the dataset, implying potential public health considerations in areas with elevated smoking and alcohol consumption rates.



**Figure3.  Adult Smoking and Alcohol Usage Distribution**

**(7) Age Distribution**

This chart shows a right-skewed distribution, it provides a clear view of the patient population, with the majority of individuals falling between 60 and 80 years old. The peak around age 70 highlights that this age group forms the largest segment of the patient base. This is consistent with the typical age profile of healthcare services that cater to older adults, such as those covered under Medicare or other elderly care programs. The distribution gradually decreases after age 80, showing fewer patients in the 90+ range. There are very few patients below the age of 50, which further emphasizes that this dataset predominantly represents an older, aging population. This insight is crucial for understanding the healthcare needs, treatment plans, and preventive care strategies that are most relevant to this elderly patient group.



**Figure4. Age Distribution**

### 3.2.2 Conclusion

These visualizations play a crucial role in understanding our dataset and form the basis for further analysis in this research.

## 3.3 Data Cleaning

In the initial phase of data cleaning, we identified relatively irrelevant variables with the following criteria and removed them from the datasets.

(1) Variables that are duplicates of other variables.
(2) Variables with low variance.
(3) Variables that are not related to our analysis.

(4) Variables that has too many Null values (the percentage of missing value > 50%)

This process helped us clean the dataset, leaving only the important and relevant data for further examination.

### 3.3.1 Handling Categorical Features and One Hot Encoding

In the dataset, member condition data was provided as categorical variables. We transformed these variables using one hot encoding. Figure following shows a graphic representation of the process.

| id | chronic_num | cond_key | cond_desc |
|---|---|---|---|
| 46388 | 8 | 226, 23, 238, 329, 48, 126, 127, 108 | Heart Failure, Except End-Stage and Acute. Other Significant Endocrine and Metabolic Disorders, Specified Heart Arrhythmias, Chronic Kidney Disease, Moderate Stage 3, Coagulation Defects and Other Specified Hematological Disorders, Dementia, Moderate, Vascular Disease |
| 172308 | 6 | 64, 38, 63, 238, 48, 226 | Cirrhosis of Liver, Diabetes with Glycemic, Unspecified, or No Complications, Chronic Liver Failure/End-Stage Liver Disorders, Specified Heart Anhythmias, Morbid Obesity |
| 1557030 | 3 | 18, 454, 17 | Cancer Metastatic to Bone, Stem Cell, Including Bone Marrow, Cancer Metastatic to Lung, |
| 1174622 | 2 | 238, 40 | Specified Heart Ambythmias, Rheumatoid Arthritis and Inflammatory Connective Tissue Disease |

**Figure5. Encoding Process**

## 3.4 Data Transformation & Feature Engineering

After removing some columns, we prepared the data for analysis by transforming it as needed, which involved the following steps:

| Web Activity | Pharmacy Usage Rate | Cost & Claim Trends | Financial Metrics |
|---|---|---|---|

| -Active Web User Indicator -Days Since Last Login | -High Pharmacy Cost Indicator -Recent Pharmacy Activity | -Cost Trend -Claims Trend -Days Since Last Claim -Total Claims | -Out-of-Pocket Cost Ratio |
|---|---|---|---|

**Figure6.  Feature Engineering**

## 1. Handling web activity data

We aimed to distinguish between active and inactive web users by creating a binary feature *active_web_user*. This feature was generated by checking if the user had any web logins (*login_pmpm_ct*). If there was at least one login, the user was considered active (*active_web_user* = 1), otherwise inactive (*active_web_user* = 0).

## 2. Creating a variable for time since last login

To address web engagement recency, we generated the feature *days_since_last_login*. Since users with no recent logins could skew the data, we assigned a high placeholder value (999) to instances where the login information was missing or invalid. This allowed us to handle the missing or negative values without removing them from the dataset.

## 3. Pharmacy cost utilization

We created the binary feature *high_pharmacy_cost* to capture high pharmacy utilization based on the distribution of *rx_overall_pmpm_cost*. Specifically, we defined high pharmacy utilization as having a cost above the median, marking these cases as *high_pharmacy_cost* = 1, while costs below or equal to the median were assigned 0.

## 4. Capturing recent pharmacy activity

For a more granular understanding of pharmacy engagement, we developed the feature *recent_pharmacy_activity*. This feature identifies patients who had a prescription in the last 30 days (*rx_days_since_last_script* <= 30), marking them as 1. All others were marked as 0. This variable served as a key indicator of recent pharmacy utilization.

## 5. Rate of change in costs and claims

We introduced two trend-based features, *cost_trend* and *claims_trend*, to account for changes over time in the patient's pharmacy costs and non-participating claims, respectively.
*cost_trend* captures the rate of change in *total_net_paid_pmpm_cost* (i.e., the percentage change from the previous time period).
*claims_trend* reflects the rate of change in non-participating claims (*nonpar_clm_ct_pmpm*). Both variables were filled with zeros when no prior data existed.

## 6. Long gaps between claims

The feature *days_since_last_clm* (already present in the dataset) was used to account for long gaps between a patient's claims. This feature helps in identifying periods of inactivity or lack of claims, which could be significant for certain analyses.

## 7. Seasonality in claim submission

While exact dates were not available, we created an approximation for seasonality by generating a categorical feature representing "low engagement periods." This feature was inferred based on

whether *total costs* or *claim counts* fell below the median during the analysis period. Additionally, we recalculated the total claims by summing up *nonpar_clm_ct_pmpm* and *oontwk_clm_ct_pmpm* to provide a comprehensive picture of claim activity through the total_claims variable.

**8. Out-of-pocket cost ratio**

We calculated the *out_of_pocket_ratio* to provide insight into how much of the total cost was borne by the patient. This ratio was derived by dividing the patient's responsibility (*total_mbr_resp_pmpm_cost*) by the total net paid costs (*total_net_paid_pmpm_cost*). This feature highlights the financial burden on patients relative to the total cost.

Finally, we have a dataset containing 381858 rows $\times$ 256 columns.

## 3.5 Feature Importance & Selection

We use random forests to generate features and assess their importance due to their ability to automatically rank and select the most relevant features in a dataset.
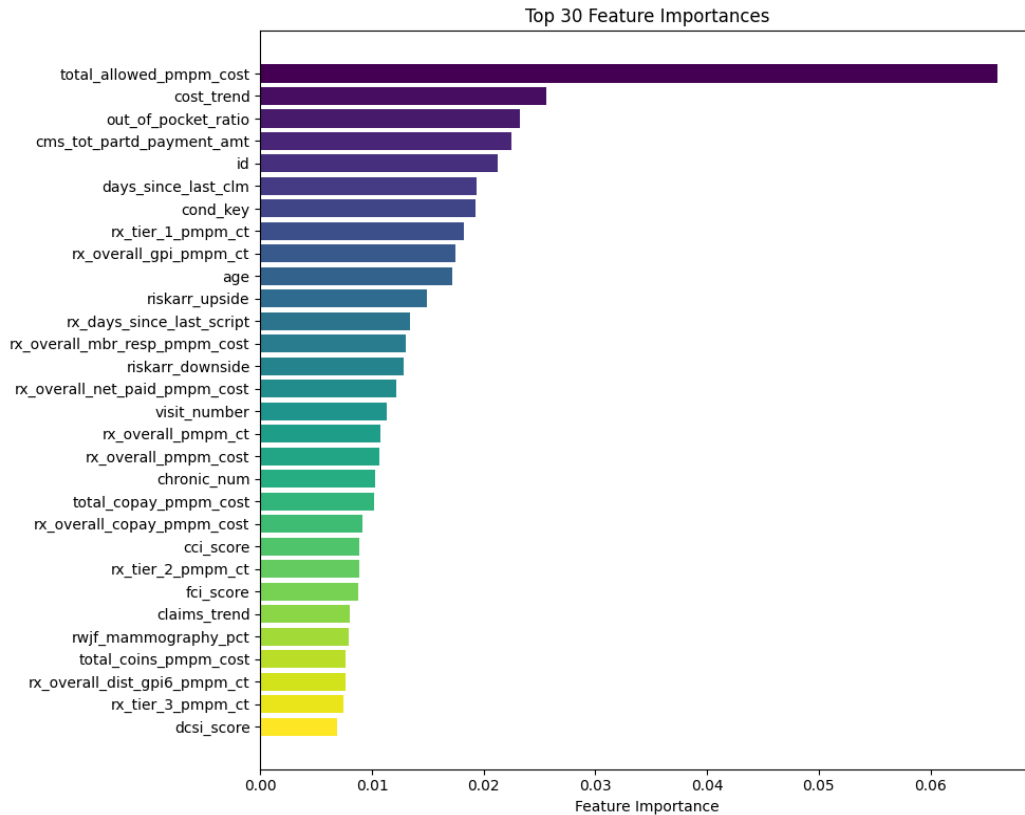
### 3.5.1 Feature Importance- Random Forests

Random forests were used to assess the importance of features in the dataset. This model ranks the features by evaluating their contribution to reducing impurity across multiple decision trees. The features that led to the greatest reduction in Gini impurity (or entropy) were deemed more important.

### 3.5.2 Advantages of Random Forests

    **(1) Feature Importance Scores**: Random forests provide a built-in mechanism for ranking feature importance. The importance scores reflect the cumulative contribution of each feature to improving the decision-making process across all trees in the forest.

    **(2) Robustness**: By averaging across multiple trees, random forests mitigate overfitting, which ensures that the feature importance results are reliable and applicable to unseen data.

    **(3) Automatic Feature Selection**: The algorithm inherently selects the most relevant features while building trees, allowing us to focus on a subset of critical features for subsequent analysis or model development.

### 3.5.3 Feature Selection Process

Based on the feature importance scores, the top features contributing significantly to the model's predictive power were selected. These features were used for further modeling and analysis, while less important features were dropped to simplify the model and improve generalization.

**Figure7. Top 30 Feature Importances**

The analysis indicates that cost-related features, such as PMPM costs and cost trends, are the primary drivers in the model. These are followed by patient risk scores, utilization patterns, and pharmacy-related features, which collectively offer a comprehensive view of patient behavior, risk, and associated costs. The model's reliance on these features highlights the importance of financial and healthcare utilization data in predicting outcomes.

### 3.5.4 Top Features Analysis

The feature importance chart figure highlights some standout features that are important. When we look at all the variables, these features can be grouped into the following categories. This grouping helps us understand how these important aspects affect the model's predictions and why they matter in our analysis. Recognizing these influential factors helps us make better decisions and improve our strategies for more accurate results. The following are the four categories of importance groups.

**1. Expenditure**: *total_allowed_pmpm_cost:* This feature dominates the importance ranking by a significant margin, indicating it has the greatest influence on the model's

predictive performance. This suggests that per-member-per-month costs (PMPM) are a strong determinant in the target variable.

*cost_trend:* The second most important feature is the cost trend, reflecting how changes in costs over time significantly affect predictions, likely indicating that historical cost increases are key predictors of future costs or outcomes.

*out_of_pocket_ratio* and *cms_tot_partd_payment_amt*: These financial-related features further emphasize the central role of cost and payment data in the analysis. The ratio of out-of-pocket expenses and CMS Part D payments (likely related to Medicare drug coverage) are highly important factors.

**2, Health Utilization Metrics**: Features such as *days_since_last_clm*, *rx_days_since_last_script,* and *visit_number* highlight the importance of patient utilization patterns, specifically how recent healthcare claims, prescription fills, and overall visit frequency contribute to the prediction.

**3, Risk Adjustment and Condition Scores:** Features such as *riskarr_upside*, *riskarr_downside,* and *cci_score* show that risk adjustment factors and health conditions are key in explaining variations in cost or outcome. These scores provide a quantitative assessment of a patient's risk based on their comorbidities or expected healthcare costs.

**4, Prescription Drug Metrics**: Several prescription-related features, such as *rx_tier_1_pmpm_ct*, *rx_overall_gpi_pmpm_ct*, and *rx_overall_copay_pmpm_cost,* highlight the importance of pharmacy data. These features indicate the number of prescriptions filled per month, categorized by drug tier, and associated costs, which are crucial in predicting outcomes.

# 4. Statistical Analysis and Modeling

## 4.1 Model Selection

To predict the target variable, '*preventive_visit_gap_ind*', a variety of advanced machine learning models were employed to conduct a comprehensive assessment of potential predictors. The dataset was split into training and validation sets, with 80% allocated for training and 20% reserved for evaluation. This method ensured that the models had sufficient data to learn effectively while maintaining an independent subset for evaluating their predictive performance.

We initially experimented with five different types of models— K-Nearest Neighbors (KNN), Logistic Regression, XGBoost, and LightGBM. Each model's performance was assessed primarily based on the Area Under the ROC Curve (AUC), a metric commonly used for binary classification problems. The AUC for the ROC curve of each model, along with a summary of the advantages and disadvantages of each method, is presented in Figure 4.1.

Before tuning the hyperparameters of each model, the LightGBM model performed the highest AUC score. XGBoost came in second, followed by KNN and Logistic Regression. In addition to having better predictive performance, the LightGBM model also demonstrated faster training times. Therefore, given its performance and efficiency, we selected LightGBM as the final model for predicting the target variable, *'preventive_visit_gap_ind'*.

| Models | AUC Score | Advantages | Disadvantages |
|---|---|---|---|
| LightGBM | 0.751 | - High efficiency in training speed<br>- Effective for distributed training<br>- Handles missing values well | - Difficult to optimize several hyperparameters<br>- Prone to overfitting<br>- Can select features with high correlation instead of the most important ones |
| XGBoost | 0.749 | - Fast execution speed<br>- Better model performance<br>- Includes regularization techniques to prevent overfitting<br>- Can handle sparse data efficiently | - Memory-intensive for large datasets<br>- Complex algorithm and hyperparameter tuning<br>- Potential for overfitting<br>- Lack of transparency in model predictions |
| Logistic Regression | 0.642 | - Simple to implement, interpret, and train efficiently<br>- Versatile, can handle both continuous and categorical variables | - May struggle with imbalanced data<br>- Requires regularization to avoid overfitting |
| KNN | 0.573 | - Simple algorithm<br>- Handles complex, non-linear relationships well | - Struggles with too many features<br>- No interpretable summary or model from training |

**Figure 7. Comparison of AUC Scores, Advantages, and Disadvantages of Models**

## 4.2 Hyperparameter Optimization- LightGBM

In order to improve the performance of our LightGBM model, we employed two different optimization techniques: Randomized Search Optimization and Bayesian Optimization. The objective was to determine the best set of hyperparameters that would maximize the AUC score on the validation set.

We applied two different optimization techniques to tune the hyperparameters of our model: Random Search and Bayesian Optimization using Optuna. Randomized Search allowed us to quickly explore a wide range of parameter combinations by randomly sampling from a defined search space, and it provided us with some promising results. In parallel, we implemented Bayesian Optimization, which uses prior knowledge from earlier trials to guide the search process more efficiently by balancing exploration and exploitation of the search space.

Both methods were evaluated by comparing the AUC scores of the models produced. After running both approaches, we found that Bayesian Optimization achieved the best result with an AUC score of 0.7536. Table # provides a comparison of the hyperparameters before and after tuning.

| Hyperparameters | Original LightGBM | RandomizedSearch Optimization | Bayesian Optimization |
|---|---|---|---|
| 'objective' | 'binary' | 'binary' | 'binary' |
| 'metric' | 'auc' | 'auc' | 'auc' |
| 'boosting_type' | 'gbdt' | 'gbdt' | 'gbdt' |
| 'learning_rate' | 0.01 | 0.01 | 0.013 |
| 'num_leaves' | 128 | 150 | 200 |
| 'min_data_in_leaf' | 20 | 30 | 25 |
| 'feature_fraction' | 0.7 | 0.5 | 0.631 |
| 'bagging_fraction' | 0.9 | 0.75 | 0.889 |
| 'bagging_freq' | 5 | 3 | 4 |
| 'lambda_l1' | 0.1 | 0.05 | 0.167 |
| 'lambda_l2' | 0.1 | 0.1 | 0.165 |

| 'scale_pos_weight' | 1 | 1.5 | None |
|---|---|---|---|
| 'max_bin' | 255 | 200 | 229 |
| 'max_depth' | 12 | 10 | 13 |
| 'min_gain_to_split' | 0.1 | 0.1 | 0.126 |
| 'early_stopping_rounds' | 50 | 30 | None |

**Notes: *objective:** Defines the loss function to optimize, such as regression or classification; **\*metric:** Specifies evaluation metrics for model performance, like accuracy or AUC; **\*boosting_type:** Selects the boosting algorithm, e.g., 'gbdt', 'dart', or 'goss'; **\*learning_rate:** Sets the step size for minimizing the loss function, affecting accuracy and training rounds; **\*num_leaves:** Maximum number of leaves in a tree, influencing model complexity and potential overfitting; **\*min_data_in_leaf:** Minimum observations required in a leaf to prevent overfitting; **\*feature_fraction:** Fraction of features to consider when building each tree, reducing overfitting; **\*bagging_fraction:** Fraction of training data to use per iteration to decrease overfitting; **\*bagging_freq:** Frequency of bagging, e.g., every 5 iterations; **\*lambda_l1:** L1 regularization to reduce overfitting with a penalty on coefficients; **\*lambda_l2:** L2 regularization with a penalty on the square of coefficients, also reducing overfitting; **\*scale_pos_weight:** Adjusts class balance in imbalanced datasets to improve sensitivity; **\*max_bin:** Maximum number of bins for feature bucketing during training; **\*max_depth:** Limits tree depth to control model complexity; **\*min_gain_to_split:** Minimum gain required for a split; below this, no split occurs; **\*early_stopping_rounds:** Rounds to wait for validation metric improvement before stopping training.
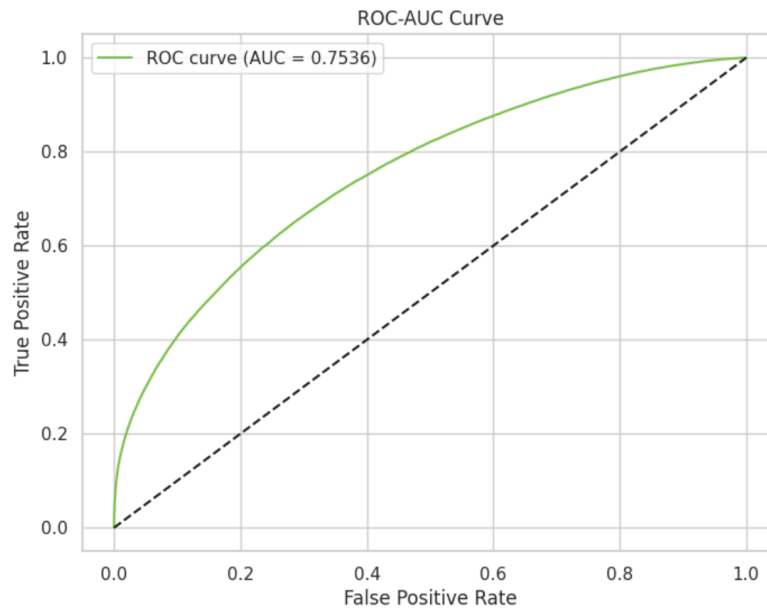
**Figure 8. Comparison of Hyperparameters Before and After Optimization**

## 4.3 Final Model Results

To gain a comprehensive understanding of our LightGBM model's performance, we've visualized the ROC-AUC curve, Precision-Recall curve, and Confusion Matrix.
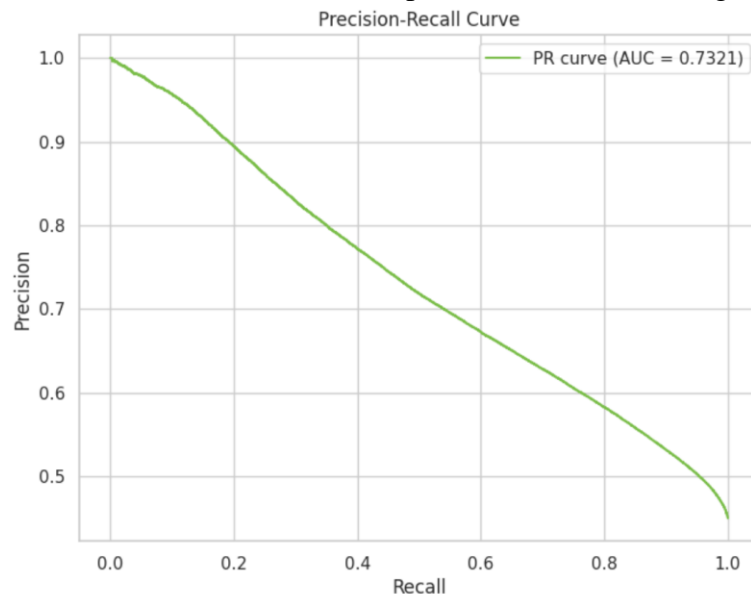
### 4.3.1 ROC-AUC Curve

The ROC-AUC curve demonstrates the model's ability to distinguish between the positive and negative classes. The AUC score of 0.7536 indicates that the model performs fairly well, with a reasonable level of discrimination between the two classes. A perfect model would have an AUC of 1, while a score of 0.5 would suggest the model is no better than random guessing. The AUC of 0.7536 suggests that the model is able to differentiate between true positives and false positives with moderate success, though there is room for improvement to make this performance stronger.

**Figure 9. ROC-AUC Curve**

### 4.3.2 Precision-Recall Curve

The Precision-Recall curve highlights the trade-off between precision and recall in the model's predictions. With a PR-AUC of 0.7321, the model demonstrates a good balance between precision (the accuracy of positive predictions) and recall (the ability to identify all actual positive cases). However, as recall increases, precision decreases, which means the model tends to lose accuracy in its positive predictions as it tries to capture more positive instances. Despite this, the overall PR-AUC score still reflects solid performance in balancing the two metrics.



**Figure 10. Precision-Recall Curve**

### 4.3.3 Confusion Matrix

The confusion matrix provides further insight into the model's predictive accuracy. The matrix shows that the model correctly identified 134,612 true negatives and 75,984 true positives, meaning that it successfully classified a large number of both positive and negative cases. However, it also made 33,487 false positive predictions (incorrectly classifying negatives as positives) and 61,396 false negative predictions (incorrectly classifying positives as negatives). The relatively high number of false negatives suggests that the model is missing a considerable number of actual positive cases, which could impact its overall usefulness, especially in cases where identifying positive cases is crucial.
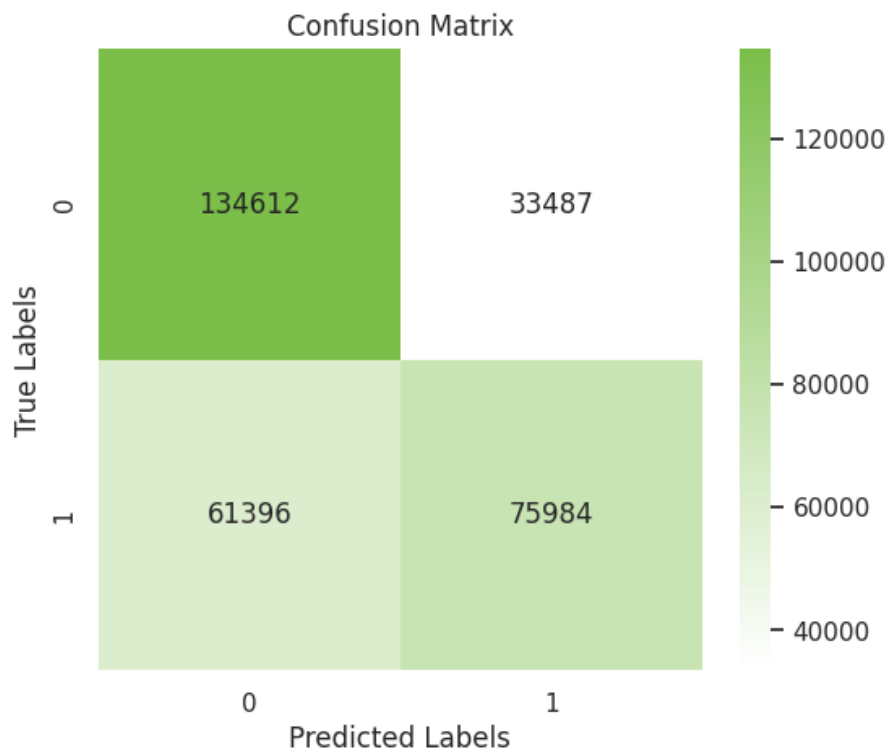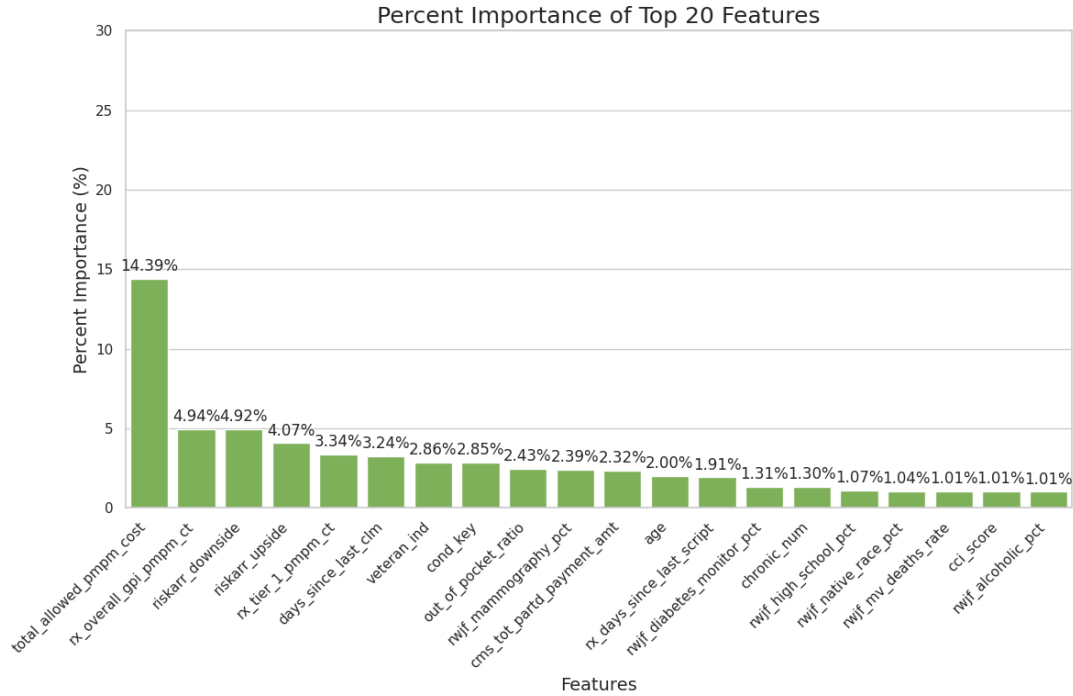


**Figure 11. Confusion Matrix**

## 4.4 Feature Importance

The feature importance plot reveals the top features influencing the model's predictions. The most significant feature was *"total_allowed_pmpm_cost"*, accounting for 14.39% of the overall importance, indicating a strong correlation with the outcomes. Following this, "rx_overall_gpi_pmpm_ct", *"riskarr_downside"* and *"riskarr_upside"* were also prominent, contributing 4.94%, 4.92%, 4.07%, respectively, reflecting their critical role in the model. Additionally, *"days_sincs_last_clm"* and *"veteran_ind"* features emerged as significant variables, with importance values of 3.34% and 3.24%, respectively. This analysis shows that cost-related features, particularly total allowed costs, are the most influential in the model's predictions.

**Figure 12. Percent Importance of Top 20 Features**

## 4.5 Limitation of Model

The current LightGBM model effectively predicts the likelihood of a member experiencing a preventive visit gap, but it does not provide insights into the timing of these gaps. Understanding when a member is likely to experience a gap would enable more proactive interventions and better resource allocation. Additionally, the model may not fully capture the complexities of patient behavior and other contextual factors that could influence adherence.

Furthermore, the model relies on selected features that may not encompass all relevant variables. For instance, certain demographic and clinical characteristics could offer valuable insights into the factors contributing to preventive visit gaps but are not included in the current dataset. The absence of these variables limits the model's ability to provide a comprehensive view of the determinants of non-engagement.

## 5. Business Implications and Recommendations

Our analysis identified several challenges faced by LPPO members in engaging with preventive care. The following recommendations are designed to address these issues realistically and effectively, aligning with Humana's operational goals while offering solutions that balance cost and member needs.

## 5.1 Refined 8-Step "Illness to Wellness" Journey

Members often disengage due to confusion about treatment, financial concerns, and logistical barriers. To address this, we propose a refined 8-step "Illness to Wellness" journey, strategically intervening where members are most likely to disengage.

### Step1: Case Manager Introduction & Disease Education

- **Practical Approach:** Assign a case manager to each LPPO member at the start of their care journey. The case manager explains the importance of preventive care, the role of the PCP, and how their health plan works.
- **Solved Problem:** Members who are high-risk for disengagement are flagged early based on their engagement patterns and socio-economic data, ensuring personalized attention from day one.

### Step2: Support Groups & Comprehensive Education

- **Practical Approach**: Connect members with peer support groups to provide emotional support. Simultaneously, offer educational resources that cover treatment plans, financial guidance, and lifestyle adjustments.
- **Solved Problem:** Peer-led groups are effective for emotional support, while professional resources focus on key medical or financial concerns. This dual approach ensures members feel supported throughout their journey.

### Step3: Transparent Treatment Costs & Options

- **Practical Approach**: Provide members with a clear breakdown of treatment costs, including co-pays, deductibles, and potential out-of-pocket expenses. Early financial transparency helps prevent surprise costs later in the journey.
- **Solved Problem:** Personalized cost breakdowns during initial consultations ensure members have a realistic understanding of their financial responsibilities, reducing the likelihood of disengagement due to cost.

### Detour 1: Addressing Concerns Early

- **Practical Approach**: If a member expresses financial, medical, or logistical concerns, Humana steps in with tailored interventions. These include offering financial assistance, providing virtual consultations, or ensuring home care support.

- **Solved Problem:** Case managers proactively identify members likely to face financial barriers by reviewing income, claims history, and past behavior, allowing early intervention before these concerns become disengagement drivers.

**Step 4: Insurance & Financial Checkpoint**

- **Practical Approach**: Case managers conduct a detailed financial review to ensure members understand their insurance coverage, co-payments, and deductibles.
- **Solved Problem:** Satisfaction scores from these financial reviews help gauge member understanding, while reduced dropout rates due to financial strain act as key performance metrics.

**Step 5: Treatment Start**

- **Practical Approach**: Members begin their treatment plan with continued case manager support. The focus is on ensuring they attend scheduled preventive visits and follow the care plan.
- **Solved Problem:**Automated check-ins and follow-ups ensure regular touchpoints without overwhelming members, helping them stay on track without feeling pressured.

**Step 6: Bi-weekly Check-ins**

- **Practical Approach**: Every two weeks, case managers check in with members to assess their progress. These check-ins serve to address any emerging concerns early and to provide reassurance about the treatment plan.
- **Solved Problem:** Using predictive analytics, Humana can flag members who miss appointments or show reduced engagement. This allows case managers to intervene quickly before issues escalate.

**Detour 2: Post-Treatment Concerns**

- **Practical Approach**: When members express concerns during treatment, case managers use personalized, data-driven tools (such as Generative AI) to provide real-time solutions, whether financial advice, medical re-consultation, or emotional support.
- **Solved Problem:** AI-powered tools allow case managers to pull detailed reports on each member's health and financial situation, enabling highly tailored interventions.

**Step 7: Positive Progress Monitoring**

- ○ **Practical Approach**: As members progress positively, case managers continue to check in, ensuring ongoing engagement with preventive care and addressing any challenges before they lead to disengagement.
- ○ **Solved Problem:** Continued emotional and informational support, even during periods of improvement, ensures members remain motivated to stay engaged in their care plan.

**Step 8: Wellness & Community Engagement**

- ○ **Practical Approach:** After recovery, members are encouraged to stay connected with the community by sharing their experiences in support groups, participating in wellness events, or mentoring others.
- ○ **Solved Problem:** This approach fosters long-term engagement even after treatment ends, creating a supportive environment that encourages members to continue preventive care and contribute to the wellness community.

## 5.2 Financial Assistance Programs

Financial concerns are a significant cause of disengagement, particularly among lower-income members. Early intervention can prevent financial strain from causing members to drop out of preventive care. Therefore, we concluded four recommendations to deal with the financial concerns.

1. **Reducing Deductibles for Initial Therapy**
   - ○ **Feasible Approach**: Offer a temporary 25% reduction in deductibles for the first six months of therapy to members identified as financially vulnerable. This helps alleviate financial pressure at the beginning of treatment without overburdening Humana.
   - ○ **Solved Problem:** This approach targets members who are most at risk of financial strain, ensuring they stay engaged without causing unsustainable financial burdens on Humana.
2. **Partial Co-insurance Waivers**
   - ○ **Feasible Approach**: Provide a 50% waiver on co-insurance for six months, focusing on essential drugs and preventive care treatments. Exclude high-cost medications like Tagrisso to keep the program financially manageable.
   - ○ **Solved Problem:** This waiver focuses on lower-cost treatments that have a significant impact on preventive care outcomes, reducing immediate financial strain for low-income members.
3. **Streamlined Claim Submission Support**

- **Feasible Approach**: Improve the claims submission process by enhancing digital platforms with real-time feedback and offering additional customer support for claim-related questions.
- **Solved Problem:** By guiding members through the claims process with better tools and support, Humana can reduce errors and rejections, keeping members engaged without financial frustration.

## 5.3 Medical Re-Consultation Support

Many vulnerable members, particularly seniors and those with disabilities, need access to flexible re-consultations and care adjustments to stay engaged with preventive care. Therefore, we concluded two recommendations to address the needs of those vulnerable members.

1. **Home-Based Care Expansion**
    - **Feasible Approach**: Expand Humana's home care services to include more preventive visits for members with mobility challenges or chronic conditions. This provides in-home consultations, reducing the need for hospital visits.
    - **Solved Problem:** Home care services will be prioritized for high-risk members, ensuring they receive necessary care without logistical challenges or stress.
2. **Virtual Care Services**
    - **Feasible Approach**: Strengthen virtual care offerings, providing 24/7 telehealth access for members requiring re-consultations, especially those dealing with side effects or needing follow-up care.
    - **Solved Problem:** Virtual care reduces unnecessary hospital visits, ensuring timely medical attention at a lower cost while keeping members engaged in their care plan.

## 5.4 Strengthening Patient Support Systems

Emotional and informational support is essential for maintaining engagement, especially for members facing health challenges or financial concerns. Therefore, we concluded two recommendations to address the needs of members with health or financial concerns.

1. **Case Manager Integration**
    - **Feasible Approach**: Assi dedicated case managers to each LPPO member. These managers use AI-driven tools to track and respond to each member's unique needs, ensuring timely interventions.
    - **Solved Problem:** AI-powered tools provide case managers with personalized insights, allowing for more effective support tailored to each member's situation.
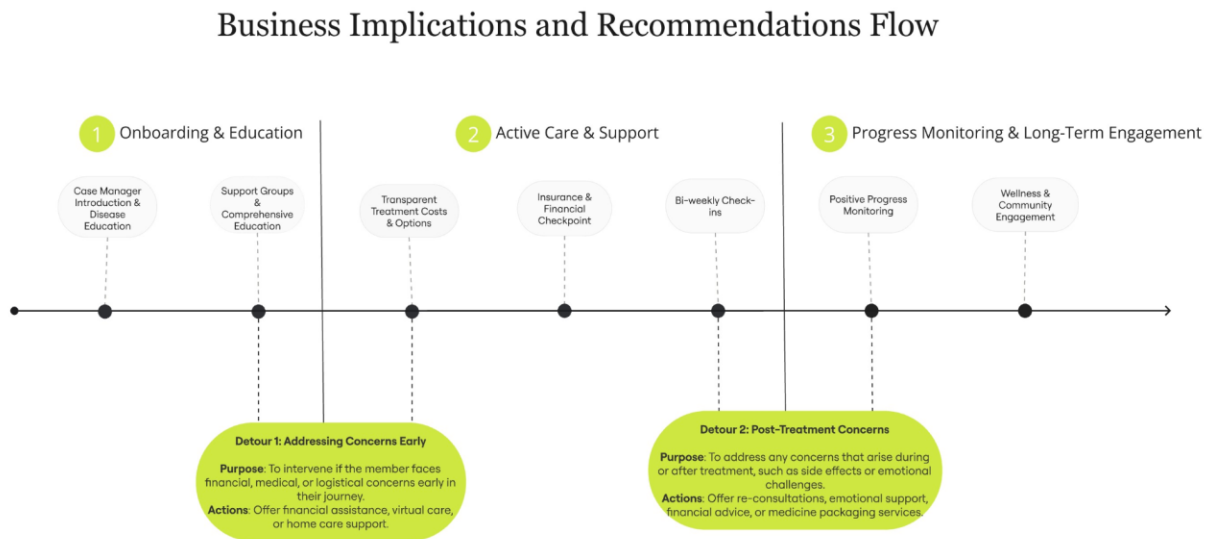2. **Structured Support Groups**
    - **Feasible Approach**: Establish structured, moderated support groups for members in preventive care, focusing on shared experiences to reduce isolation and offer emotional reassurance.

○ **Solved Problem:** These groups offer a mix of peer support and professional guidance, helping members stay motivated and engaged with their care plans.

## 5.5 Conclusions

By focusing on early financial assistance, expanding medical re-consultation options, and strengthening patient support systems, Humana can significantly improve LPPO member engagement with preventive care. These realistic and actionable recommendations will lead to better health outcomes, cost savings, and higher member satisfaction, positioning Humana for success in the Medicare Advantage market.



**Figure13. Business Implications and Recommendations Flow**

## 6. Conclusions

This report highlights the significant challenges Humana's LPPO members face in staying engaged with preventive care. By focusing on practical solutions such as improving the member experience, providing financial support, and offering stronger medical guidance, our

recommendations aim to increase engagement rates, ultimately improving health outcomes and ensuring financial sustainability for Humana.

Our suggested strategies are designed to help members navigate the often confusing landscape of healthcare by addressing key pain points like financial strain, lack of clarity in treatment plans, and emotional support gaps. The ultimate goal is to keep more members on track with preventive care, which not only benefits their health but also improves Humana's Stars ratings and overall performance.

If implemented effectively, these recommendations will enhance the member experience and reduce the risk of disengagement. The result will be a stronger, more sustainable Medicare Advantage program that delivers both improved member satisfaction and business success for Humana.

## 7. References

1. MyEducator. (n.d.). *Advantages and disadvantages of KNN*. MyEducator. Retrieved October 20, 2024, from https://app.myeducator.com/reader/web/1421a/11/q07a0/
2. Brownlee, J. (2016, May 9). *A gentle introduction to XGBoost for applied machine learning*. Machine Learning Mastery. Retrieved October 20, 2024, from

https://machinelearningmastery.com/gentle-introduction-xgboost-applied-machine-learning/

3. Simplilearn. (2023, July 10). *What is XGBoost algorithm in machine learning?* Simplilearn. Retrieved October 20, 2024, from https://www.simplilearn.com/what-is-xgboost-algorithm-in-machine-learning-article

4. Krayonnz. (2022, March 23). *What are the advantages and disadvantages of XGBoost?* Krayonnz. Retrieved October 20, 2024, from https://www.krayonnz.com/user/doubts/detail/623b2b7235e21e005f953106/what-are-the-advantages-and-disadvantages-of-XGBoost

5. GeeksforGeeks. (2021, November 1). *Advantages and disadvantages of logistic regression*. GeeksforGeeks. Retrieved October 20, 2024, from https://www.geeksforgeeks.org/advantages-and-disadvantages-of-logistic-regression/

6. Kaggle. (2022, September 15). *Advantages and disadvantages of logistic regression*. Kaggle Discussions. Retrieved October 20, 2024, from https://www.kaggle.com/discussions/general/352871

7. Kaggle. (2021, March 22). *Advantages and disadvantages of LightGBM*. Kaggle Discussions. Retrieved October 20, 2024, from https://www.kaggle.com/discussions/general/264327

8. Data Aspirant. (2023, June 23). *LightGBM algorithm: A detailed introduction*. Data Aspirant. Retrieved October 20, 2024, from https://dataaspirant.com/lightgbm-algorithm/#t-1679668681677

9. LightGBM. (n.d.). *LightGBM hyperparameters*. LightGBM Documentation. Retrieved October 20, 2024, from https://lightgbm.readthedocs.io/en/latest/Parameters-Tuning.html

10. Miro. (n.d.). *Miro board dashboard*. Miro. Retrieved October 20, 2024, from https://miro.com/app/dashboard

11. Builtin. (n.d.). *Feature engineering: Empowering data through transformation*. Builtin. Retrieved October 20, 2024, from https://builtin.com/articles/feature-engineering

12. Scikit-learn. (n.d.). *sklearn.model_selection.RandomizedSearchCV*. Scikit-learn documentation. Retrieved October 20, 2024, from https://scikit-learn.org/1.5/modules/generated/sklearn.model_selection.RandomizedSearchCV.html