

# Census Income Classification and Segmentation

JPMorgan Chase Data Science Challenge

Amir Taherkhani

November 5, 2025

## Contents

<b>1</b>	<b>Executive Summary</b>	<b>2</b>
<b>2</b>	<b>Data Understanding and Quality</b>	<b>2</b>
2.1	The Challenge of Missing Information . . . . .	2
2.2	The Class Imbalance Problem . . . . .	2
2.3	Data Preparation Steps . . . . .	3
2.4	Technical Considerations . . . . .	3
<b>3</b>	<b>Classification Model Development</b>	<b>3</b>
3.1	Model Selection and Training . . . . .	3
3.2	Understanding Model Performance . . . . .	4
3.3	The Threshold Question . . . . .	5
3.4	What Drives the Predictions . . . . .	6
3.5	Technical Considerations . . . . .	7
<b>4</b>	<b>Customer Segmentation Analysis</b>	<b>8</b>
4.1	Methodology and Segment Selection . . . . .	8
4.2	Segment Profiles . . . . .	9
4.3	Technical Considerations . . . . .	9
<b>5</b>	<b>Conclusions and Recommendations</b>	<b>10</b>
5.1	Limitations and Considerations . . . . .	10

# 1 Executive Summary

This analysis addresses a fundamental question for customer targeting: how can we identify high-income individuals and understand different customer groups? We developed two complementary approaches. The first approach uses a predictive model that identifies individuals earning at least \$50,000 with 93.8% accuracy. The second approach groups customers into three segments based on their life stages, enabling targeted marketing strategies for each group.

The key finding is that effective prediction depends not just on model accuracy, but on choosing the right decision threshold. While the model can correctly identify patterns in the data, the question of where to draw the line between high and low income predictions has significant practical implications. We found that the commonly used default threshold optimizes statistical measures but may not align with business objectives. Different thresholds favor either reaching more potential customers or focusing on the most certain cases. Thus, we analyzed three threshold options, each balancing precision and recall differently.

The customer segmentation reveals three distinct groups: older adults with reduced work participation, working professionals with stable employment, and younger individuals with limited labor market attachment. These groups show dramatically different income patterns, with high earners representing 13.3% of working professionals but virtually absent from the younger group.

## 2 Data Understanding and Quality

The analysis began with a dataset containing approximately 200,000 records from the 1994-1995 U.S. Census Bureau Current Population Surveys. Each record represents an individual with 40 different characteristics including age, education, work patterns, and financial information. The dataset also includes population weights, which account for how the survey was designed to represent the broader U.S. population rather than just the people who happened to be surveyed.

### 2.1 The Challenge of Missing Information

Working with real-world survey data meant confronting how missing or non-applicable information was recorded. The dataset used several different ways to indicate when a question did not apply to someone or when information was missing. For example, fields like "industry" and "occupation" would show "Not in universe" for people who were not employed, since questions about their workplace would not make sense. Birth country information showed "?" for 4% of records where this information was unknown. Other fields used the code 0 to indicate the question was not applicable.

This variety of missing data indicators created a challenge. If we treated "Not in universe" as just another category of response rather than as missing data, it could mislead the model. Someone who is unemployed is fundamentally different from someone employed in retail or manufacturing. The preprocessing pipeline recognized these patterns and handled them appropriately.

### 2.2 The Class Imbalance Problem

Perhaps the most significant challenge was the distribution of income levels in the data. Of all individuals in the dataset, 93.8% earned less than \$50,000 annually, while only 6.2% earned at or above this threshold. This creates a 15-to-1 imbalance ratio.

This imbalance matters because a naive model could simply predict that everyone earns less than \$50,000 and achieve 93.8% accuracy without learning anything meaningful about what distinguishes high earners. Standard accuracy metrics become misleading in this context. We need methods that specifically account for the minority class we are trying to identify.

## 2.3 Data Preparation Steps

We removed 3,229 duplicate records, representing 1.6% of the dataset. Duplicates artificially inflate model confidence by giving extra weight to identical cases.

For missing values in numeric fields like age and hours worked, we used the median value of all non-missing cases. The median is less affected by extreme values than the average, making it more reliable when dealing with skewed data. For categorical fields like education and occupation, we used the most common value. This approach preserves the dominant patterns in the data while filling gaps.

A key decision involved how to represent categorical information. Education has 17 different levels, from elementary school through doctorate degrees. Occupation has 15 categories. One common approach creates a separate binary column for each possible value, but this would generate over 200 columns for our 32 categorical features. Instead, we converted categories to numeric codes that preserve order where it exists. High school comes before bachelor's degree, which comes before master's degree. This keeps the data manageable while retaining meaningful relationships.

We split the data into training and testing sets using an 80/20 division. We maintained the 15-to-1 income ratio in both sets to ensure consistent conditions. The population weights were preserved throughout to ensure results generalize to the full population rather than just our sample.

Finally, we standardized numeric features by converting them to a common scale where the average is 0 and most values fall between -3 and +3. This prevents features measured in different units from dominating the model simply because they have larger numbers.

## 2.4 Technical Considerations

The preprocessing pipeline involved several technical decisions that balance data quality with model performance:

- **Missing Value Strategy:** Applied median imputation for continuous variables and mode imputation for categorical variables after converting implicit missing indicators ("Not in universe", "?", 0) to explicit null values.
- **Encoding Method:** Label encoding was chosen over one-hot encoding to manage dimensionality. With education having 17 levels and occupation 15 categories, one-hot encoding would create 200+ sparse columns. Label encoding preserves ordinal relationships where they exist.
- **Stratified Sampling:** The train-test split used stratification to maintain the 93.8%/6.2% income distribution in both sets, preventing evaluation bias from random sampling variations.
- **Standardization:** Applied z-score normalization to continuous features only (age, weeks worked, wages, capital gains, losses, dividends), leaving encoded categorical variables on their natural scale.

# 3 Classification Model Development

With clean data prepared, we turned to building a model that could predict whether someone earns at least \$50,000 based on their demographic and employment characteristics. The goal was not just accuracy, but understanding what patterns drive high income and how confidently we can make predictions in different cases.

## 3.1 Model Selection and Training

We selected XGBoost because it handles mixed data types well, working with both categorical features like occupation and numeric features like age without requiring extensive transformation.

The model architecture includes several important features for our specific problem. First, it directly addresses class imbalance by giving extra weight to the minority class during training. We set this weight to 14.85, matching the ratio of low to high earners in our data. This prevents the model from simply predicting everyone as low income. Second, the model incorporates sample weights from the Census Bureau, ensuring our predictions reflect population patterns rather than just our sample. Third, it optimizes for F1-score, which balances catching as many high earners as possible while maintaining reasonable precision.

Training used 20 iterations of parameter optimization to find the best configuration. The optimization used 3-fold cross-validation, meaning we split the training data into thirds and trained on two-thirds while validating on the remaining third, rotating through all combinations.

The final model configuration emerged from this process: trees can grow up to 10 levels deep, providing enough complexity to capture patterns without overfitting to noise. The learning rate of 0.076 controls how much each new tree adjusts the predictions. We use 500 trees total, finding this balances accuracy with training time. Two regularization parameters prevent overfitting by randomly sampling 72% of data points and 73% of features for each tree, ensuring trees learn different aspects of the patterns.

## 3.2 Understanding Model Performance

The model achieved strong differentiation ability, with an ROC-AUC score of 0.938. This metric measures how well the model separates high and low earners across all possible threshold settings, with 1.0 being perfect and 0.5 being random chance. The PR-AUC score of 0.630 provides a complementary view focused on performance with imbalanced classes.

However, these overall scores mask an important reality. At the default threshold of 0.5, where the model classifies anyone with a predicted probability above 50% as a high earner, precision drops to 47.8%. This means that if we contacted 100 people the model classified as high earners, only 48 would actually earn at least \$50,000. The remaining 52 would be false positives.

At the same time, recall reaches 67.5%, meaning the model successfully identifies about two-thirds of all actual high earners in the data. These two metrics pull in opposite directions. We can increase precision by being more selective, but this means missing more actual high earners. We can increase recall by being more inclusive, but this means more false positives.

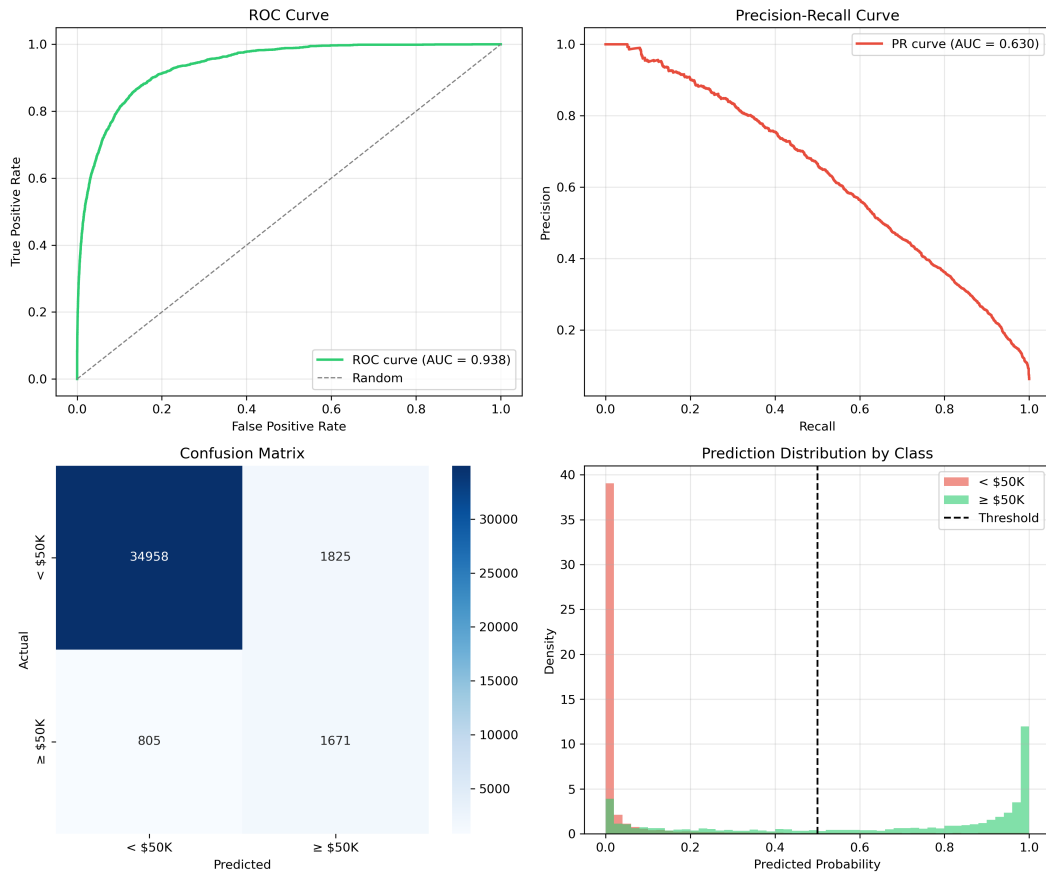


Figure 1: Model performance curves showing strong differentiation ability with ROC-AUC of 0.938.

### 3.3 The Threshold Question

This brings us to a fundamental question: where should we set the threshold? The default value of 0.5 assumes equal cost for two types of errors, but this rarely reflects real situations. In practice, the costs and benefits of correctly identifying a high earner versus incorrectly targeting a low earner are different.

We analyzed performance across different threshold values to understand the tradeoffs. At a threshold of 0.71, we achieve balanced performance with 58% precision and 58% recall. This represents the point where the F1-score, a combined measure of precision and recall, reaches its maximum.

Moving to a more selective threshold of 0.81 increases precision to 64.9% while recall decreases to 51.8%. This means we contact fewer people overall, but a higher percentage of them are actually high earners. We successfully reach about half of all high earners while maintaining nearly two-thirds accuracy in our positive predictions.

An even more conservative threshold of 0.87 pushes precision to 70.1% but drops recall to 46.4%. At this point, we are highly confident in the people we target, but we miss more than half of potential high earners.

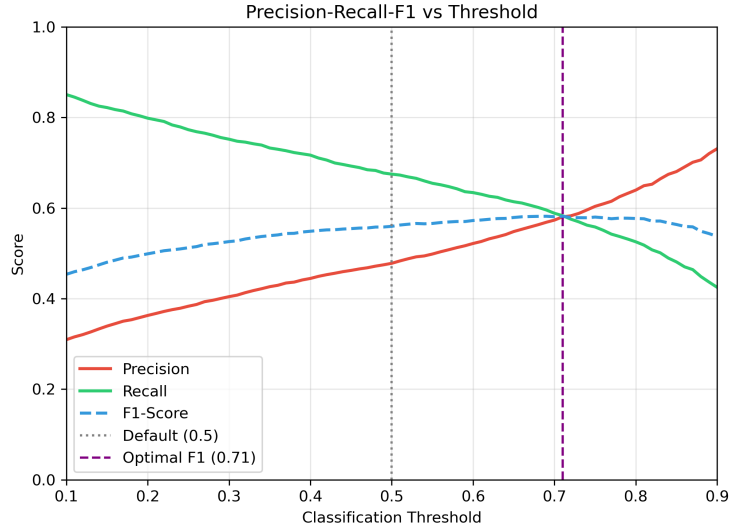


Figure 2: Performance across different thresholds. The optimal choice depends on the relative importance of precision versus recall for the specific application.

The choice between these thresholds depends on the relative importance of reaching more potential customers versus focusing resources on the most certain cases. There is no single "correct" answer, only tradeoffs that align differently with different objectives.

### 3.4 What Drives the Predictions

Understanding which features most influence the model's predictions provides insights into income patterns in the data. The number of weeks worked per year dominates at 16% of total feature importance. This makes intuitive sense as full-year employment strongly signals stable income.

Occupation ranks second at 6.7%, indicating that what you do matters more than many demographic characteristics. The type of work provides information about income level beyond what we know about age, education, or other factors.

Capital gains (4.9%) and dividend income (3.2%) serve as wealth indicators. People with investment income tend to have higher overall earnings. Sex appears fourth at 4.5%, reflecting wage gap patterns in the 1990s data.

Education appears seventh at 2.6% despite being one of the strongest single predictors of income when examined alone. This occurs because much of education's effect works through other variables already in the model. Higher education leads to professional occupations, which leads to more weeks worked and higher wages. Once these intermediate factors are included, education's direct additional contribution becomes smaller.

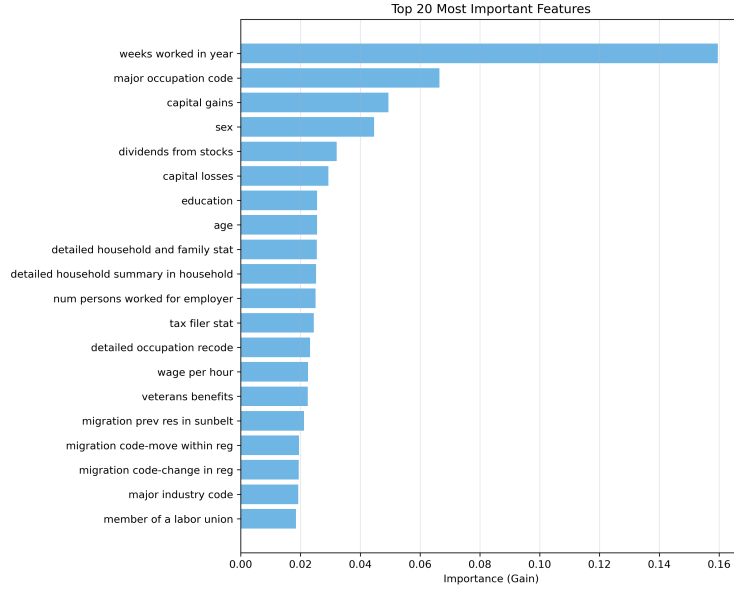


Figure 3: Features ranked by their contribution to model predictions. Work engagement and occupation dominate, while demographic factors play supporting roles.

Cross-validation showed consistent performance across different subsets of the training data, with F1-scores ranging from 0.54 to 0.58. Test set performance matched training performance, indicating the model generalizes well to new data without overfitting.

### 3.5 Technical Considerations

Several technical aspects of the modeling process deserve attention:

- **Algorithm Selection:** XGBoost was chosen over logistic regression and random forests for its superior handling of mixed feature types, built-in class imbalance weighting, and native sample weight support. The gradient boosting framework iteratively improves predictions by learning from residual errors.
- **Hyperparameter Optimization:** Used Bayesian optimization with 20 iterations and 3-fold cross-validation. This approach is more sample-efficient than grid search, particularly important with 157,000 training samples. Optimized for F1-score rather than accuracy due to class imbalance.
- **Final Parameters:** `max_depth=10` controls tree complexity, `learning_rate=0.076` sets gradient step size, `n_estimators=500` provides 500 boosting iterations, `subsample=0.72` and `colsample_bytree=0.73` provide stochastic regularization, `min_child_weight=10` prevents overfitting on minority class, `gamma=0.5` adds additional regularization.
- **Class Imbalance Handling:** `scale_pos_weight=14.85` directly addresses the 15:1 class imbalance by upweighting positive class loss during training. Without this adjustment, the model would achieve high accuracy by predicting the majority class.
- **Evaluation Metrics:** ROC-AUC measures differentiation ability across all thresholds, PR-AUC focuses on performance with imbalanced classes, F1-score balances precision and recall, and confusion matrices show actual prediction counts at specific thresholds.
- **Feature Importance:** Calculated using gain-based importance, measuring the average gain of splits using each feature across all trees. Education’s lower direct importance despite strong univariate correlation reflects its mediated effect through occupation and work patterns.

## 4 Customer Segmentation Analysis

While the classification model predicts individual income levels, segmentation takes a different approach. Instead of asking whether someone is high or low income, it asks what natural groups exist in the customer base and how these groups differ from each other. This provides a complementary view for developing marketing strategies.

### 4.1 Methodology and Segment Selection

We applied K-means clustering, a method that groups similar individuals together based on 15 characteristics. These characteristics span demographics, employment patterns, financial behavior, and household structure.

A key question was how many segments to create. We evaluated solutions with three, four, and five clusters using multiple quality metrics. While five clusters achieved slightly better statistical separation, we selected three clusters for two important reasons.

First, the five-cluster solution contained redundant groups. Two segments showed nearly identical income distributions at 10.9% and 12.4% high earners. Another pair showed 0.03% and 1.8% high earners. These similarities suggested we were splitting natural groups unnecessarily.

Second, three segments align with recognizable lifecycle stages that marketing teams can operationalize. Five segments with subtle statistical differences would be harder to distinguish and target in practice.

The silhouette score of 0.163 indicates weak statistical separation between clusters. Statistical separation is not the primary goal; what matters is whether the segments show meaningful business differences.

On this criterion, the three-segment solution succeeds. The segments show a 292-fold difference in high-income rates, from 0.05% to 13.3%. Work engagement varies systematically across segments. Investment behavior shows clear patterns. These practical differences justify the segmentation despite modest statistical separation.

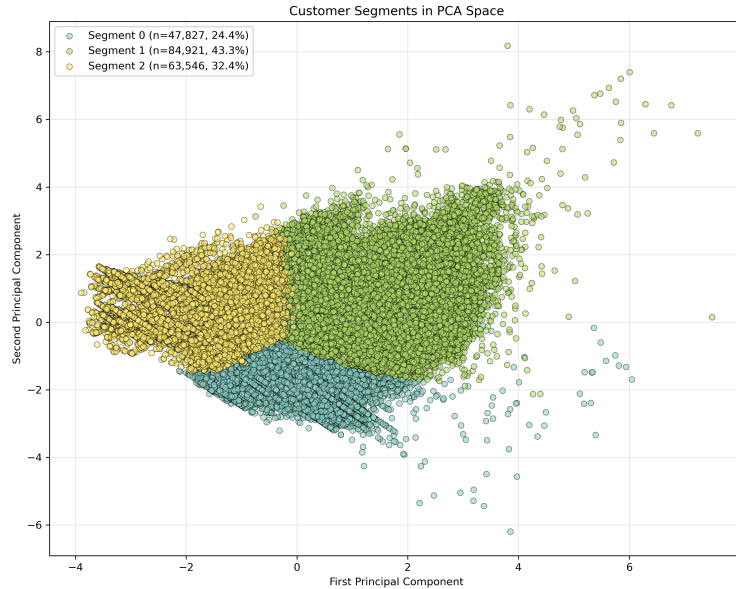


Figure 4: Visualization of customer segments using principal component analysis. The first two components explain 27.9% of variation, indicating high-dimensional complexity that cannot be fully captured in two dimensions.

The visualization uses principal component analysis to project the 15-dimensional clustering space onto two dimensions for viewing. The first component explains 18.1% of variance, the second explains 9.8%, for a total of 27.9%. The relatively low explained variance indicates the data

has high-dimensional structure that cannot be simplified to two dimensions without information loss. The segments overlap in the visualization but maintain distinct centers, reflecting that they represent tendencies rather than hard boundaries.

## 4.2 Segment Profiles

1. **Segment 0: Older Adults and Retirees (47,827 people, 24.4%)** This segment shows reduced labor market participation with accumulated financial assets. Only 2.1% earn at least \$50,000, the lowest across segments, and work engagement is 71% below average. Age averages one standard deviation above the population. Investment patterns are notable: 10.5% receive dividend income versus 3.4% with capital gains, suggesting buy-and-hold strategies typical of retirees rather than active trading. The profile fits individuals who have withdrawn from full-time employment but maintain financial resources through retirement accounts, social security, pensions, and long-term investments.
2. **Segment 1: Working Professionals (84,921 people, 43.3%)** This segment represents stable employment with the highest income potential across groups. At 13.3%, the high-income rate is six times the next highest segment. These individuals work full-year schedules with average work engagement. Investment activity is moderate: 6.5% show capital gains and 8.1% receive dividend income. Age trends slightly above average, suggesting mid-career positioning. The segment includes mid-career professionals, managers, skilled tradespeople, and business owners working full schedules and building financial assets. This large segment likely includes a premium subgroup with substantially higher income and investment activity that the three-cluster solution merges into the broader category.
3. **Segment 2: Early Career and Low Attachment (63,546 people, 32.4%)** This segment shows limited labor market attachment with minimal current income potential. Just 0.05% earn high incomes and work engagement is 81% below average. Investment activity is nearly absent: 0.3% show capital gains and 0.4% receive dividends. Age is one standard deviation below average. The combination of young age and low work engagement signals limited labor market attachment. This profile includes students, part-time workers, people between jobs, and early-career individuals without stable employment. The near-zero high-income rate indicates minimal current purchasing power, though some will transition to higher-earning segments as careers develop.

## 4.3 Technical Considerations

The segmentation methodology involved several technical decisions:

- **Algorithm:** K-means clustering was selected for its interpretability and efficiency with large datasets. The algorithm minimizes within-cluster sum of squared distances through iterative assignment and centroid updates.
- **Features:** Used 15 features capturing demographics (age, sex, race, marital status, education), employment (weeks worked, occupation, industry, class of worker), financial behavior (capital gains, losses, dividends, wage per hour), and household structure (family status, children under 18).
- **Cluster Number Selection:** Evaluated  $k=3$ , 4, and 5 using three metrics:
  - Silhouette score: 0.163 for  $k=3$  (measures cluster cohesion and separation)
  - Calinski-Harabasz index: 26,116 (higher is better, measures ratio of between-cluster to within-cluster variance)
  - Davies-Bouldin index: 2.12 (lower is better, measures average similarity between clusters)

- **Standardization:** All features were z-score normalized before clustering to prevent features with larger scales from dominating distance calculations.
- **Initialization:** K-means used k-means++ initialization to select starting centroids intelligently, reducing sensitivity to random initialization.
- **Convergence:** Algorithm ran until centroids moved less than 0.0001 or maximum 300 iterations reached.
- **Dimensionality Reduction:** PCA was used purely for visualization. Clustering was performed in the original 15-dimensional space to preserve full information. The 27.9% variance explained by two principal components indicates the high-dimensional nature of customer differences.

## 5 Conclusions and Recommendations

This analysis demonstrates that effective customer targeting requires both prediction and segmentation working together. The classification model excels at identifying high-income individuals with 93.8% differentiation ability, but its practical value depends critically on threshold selection. The segmentation model identifies natural customer groups that align with lifecycle stages, enabling targeted approaches even when individual income prediction is uncertain.

### 5.1 Limitations and Considerations

Several limitations require attention when applying these findings:

**Data Currency:** The 1994-1995 source data is now three decades old. Income distributions, wage gaps, education returns, and financial behaviors have evolved substantially. The gender wage gap has narrowed. Technology has changed how people work and invest. The model patterns may not fully reflect current conditions.

**Protected Characteristics:** The model uses sex and race as predictive features. While these improve predictions, their inclusion creates potential regulatory concerns. Deployment would require careful legal review and consideration of fair lending and discrimination laws. Alternative approaches might exclude these features or apply fairness constraints.

**Assumption Sensitivity:** The threshold analysis depends on assumed costs and conversion rates. If actual values differ significantly, the optimal threshold shifts. Organizations should validate assumptions through pilot testing before full deployment.

**Population Generalization:** Sample weights ensure the model reflects the survey population, but the Census sample may not perfectly represent specific customer bases. Regional, industry, or demographic differences between the model population and actual customers could affect performance.