

UDACITY DATA SCIENCE NANO DEGREE CAPSTONE PROJECT.

STARBUCKs

Intro

This project aims to enhance sales efficiency for Starbucks by analyzing customer behavior related to offers sent through its rewards mobile app. By examining transaction, demographic, and offer data, we seek to identify patterns that predict offer completions and determine which demographic groups respond best to various types of offers. Understanding these patterns will enable us to tailor marketing strategies more effectively, thereby boosting revenue and optimizing the impact of promotional efforts. Addressing this problem is significant as it can lead to more personalized customer interactions, higher engagement, and increased sales for Starbucks.

Description of Input Data

We’re provided with a simulated data set from the Starbucks rewards mobile app, stored in three JSON files: `portfolio.json`, `profile.json`, and `transcript.json`. These files contain information about offers, customer demographics, and transaction records, respectively.

The `portfolio.json` file details various offers, including their type (e.g., BOGO, discount, informational), difficulty, reward, duration, and channels. This information helps us understand the nature and requirements of each offer.

portfolio.head()

	reward	channels	difficulty	duration	offer_type	id
0	10	[email, mobile, social]	10	7	bogo	ae264e3637204a6fb9bb56bc8210ddfd
1	10	[web, email, mobile, social]	10	5	bogo	4d5c57ea9a6940dd891ad53e9dbe8da0
2	0	[web, email, mobile]	0	4	informational	3f207df678b143eea3cee63160fa8bed
3	5	[web, email, mobile]	5	7	bogo	9b98b8c7a33c4b65b9aebf6a799e6d9
4	5	[web, email]	20	10	discount	0b1e1539f2cc45b7b9fa7c272da2e1d7

The `profile.json` file includes demographic data for each customer, such as age, gender, income, and the date they became a member. This allows for customer segmentation based on demographic factors.

profile.head()

	gender	age	id	became_member_on	income
0	None	118	68be06ca386d4c31939f3a4f0e3dd783	20170212	NaN
1	F	55	0610b486422d4921ae7d2bf64640c50b	20170715	112000.0
2	None	118	38fe809add3b4fcd9315a9694bb96ff5	20180712	NaN
3	F	75	78afa995795e4d85b5d9ceeca43f5fef	20170509	100000.0
4	None	118	a03223e636434f42ac4c3df47e8bac43	20170804	NaN

The `transcript.json` file records all customer interactions with the app, including transactions and offer-related events. This data is crucial for tracking customer actions and linking them to specific offers.

transcript.head()

	person	event	value	time
0	78afa995795e4d85b5d9ceeca43f5fef	offer received	{'offer id': '9b98b8c7a33c4b65b9aebfe6a799e6d9'}	0
1	a03223e636434f42ac4c3df47e8bac43	offer received	{'offer id': '0b1e1539f2cc45b7b9fa7c272da2e1d7'}	0
2	e2127556f4f64592b11af22de27a7932	offer received	{'offer id': '2906b810c7d4411799c6938adc9daaa5'}	0
3	8ec6ce2a7e7949b1bf142def7d0e0586	offer received	{'offer id': 'fafdc668e3743c1bb461111dcafc2a4'}	0
4	68617ca6246f4fbc85e91a2a49552598	offer received	{'offer id': '4d5c57ea9a6940dd891ad53e9dbe8da0'}	0

By integrating these datasets, we analyzed customer behavior and identified patterns that predict offer completions, helping to tailor marketing strategies and increase sales efficiency for Starbucks.

Strategy for solving the problem.

Overall, the approach involved a comprehensive data-driven process, starting from data cleaning and EDA to feature engineering, model selection, hyperparameter tuning, evaluation, and visualization. Business questions were addressed throughout the process to ensure alignment with the goal, guiding decisions at each stage of analysis and modeling.

1. Data Cleaning and Preprocessing: The initial step involved cleaning the raw dataset to ensure its quality and integrity. It included handling missing values, removing duplicates, and ensuring consistency in data formats. Missing values were handled by removing rows/columns with missing data. Duplicate records were not found.

Data consistency checks were performed to ensure that categorical variables were properly labeled, numerical values were within expected ranges, and timestamps were formatted correctly. Preprocessing tasks also involve encoding categorical variables, such as offer types and customer demographics, to prepare them for analysis and modeling.

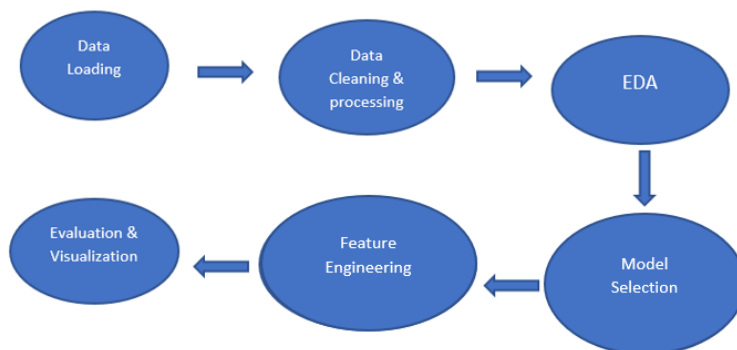
2. Exploratory Data Analysis (EDA): EDA provided insights into the dataset's structure and identified patterns and trends using summary statistics and visualization techniques. This helped address business questions like preferred offers by gender, age group impacts, effective channels for conversions, and relationships between income level and offer completion.

3. Model Selection and Hyperparameter Tuning: Once the dataset was cleaned, the appropriate modeling technique was selected based on the problem requirements, dataset characteristics, and business context. XGBoost was chosen as the baseline model due to its ability to handle complex datasets, feature importance analysis, and interpretability. Hyperparameter tuning was performed using techniques like Grid Search Cross-Validation to optimize model performance and generalization. Parameters such as `max_depth`, `learning_rate`, and `n_estimators` were tuned to improve the model's accuracy. Model selection and hyperparameter tuning were iterative processes, guided by evaluation metrics and business objectives.

4. Feature Engineering: Created new variables to capture customer behavior, offer characteristics, temporal patterns, and historical trends, enriching the dataset and enhancing predictive performance.. Feature engineering aimed to capture complex relationships and behavioral patterns that could not be adequately represented by raw data alone.

5. Evaluation and Visualization: Model performance was evaluated using accuracy, precision, recall, F1-score, confusion matrices, precision-recall curves, ROC curves, cross-validation scores, and learning curves. Visualization techniques helped interpret model behavior and communicate insights effectively.

Expected Outcome is a predictive model that accurately identifies customer responses to Starbucks' offers, informing marketing strategies, optimizing offer targeting, and enhancing customer engagement.



Metrics with justification:

The selected metrics provide insights into the model's predictive accuracy, its ability to distinguish between different classes, and its overall performance. Here are the evaluation metrics used in this context along with their justifications:

1. Accuracy: measures the proportion of correctly classified instances out of the total instances. It provides an overall assessment of the model's correctness. Accuracy gives a general idea of how well the model performs across all classes. However, it might not be sufficient when dealing with imbalanced datasets, as high accuracy can be achieved by predicting the majority class.

2. Precision and Recall: Precision measures the proportion of true positive predictions out of all positive predictions made by the model. It quantifies the model's ability to avoid false positives. Recall (or sensitivity) measures the proportion of true positive predictions out of all actual positive instances. It quantifies the model's ability to capture all positive instances. Precision reflects the accuracy of predicting offer completions, which is crucial for resource allocation and cost-effectiveness. Recall indicates the model's effectiveness in capturing all offer completions, which is essential for maximizing customer engagement and revenue.

3. F1-score: is the harmonic mean of precision and recall, providing a balance between the two metrics. It represents the model's overall performance, considering both false positives and false negatives. It accounts for both precision and recall. It is particularly useful when there is an imbalance between the classes or when both false positives and false negatives need to be minimized. In this context achieving a high F1-score ensures a balance between correctly identifying offer completions and minimizing misclassifications.

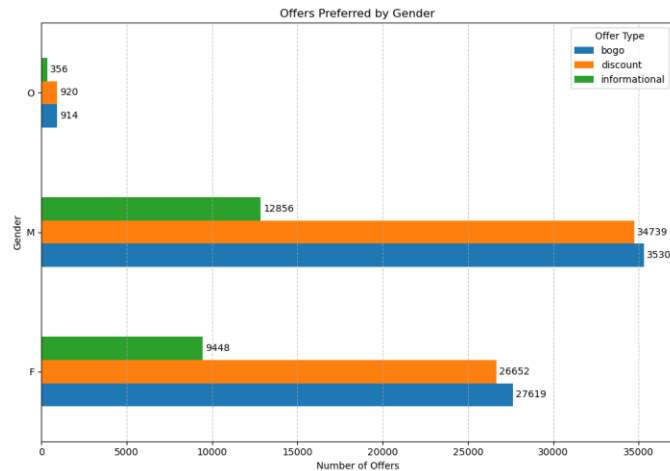
4. Area Under the ROC Curve (AUC-ROC): measures the ability of the model to distinguish between positive and negative classes across different thresholds. It quantifies the model's ability to rank instances correctly. AUC-ROC is useful for assessing the discriminatory power of the model, especially when dealing with imbalanced datasets. In the case of Starbucks datasets provided, where the number of offer completions is relatively low compared to non-completions, AUC-ROC provides insights into how well the model separates positive and negative instances.

Exploratory Data Analysis (EDA)

EDA was conducted to gain insights into the data and answer the following business questions:

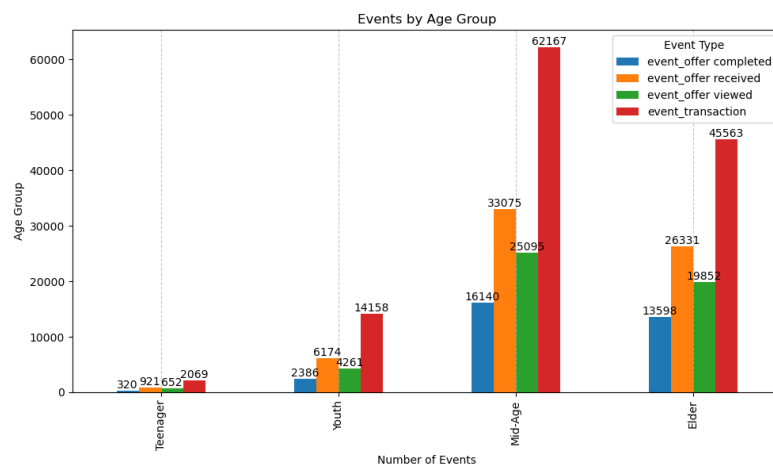
- What offers are mostly preferred by different genders?

Overall, both males and females prefer bogo and discount offers significantly more than informational offers.



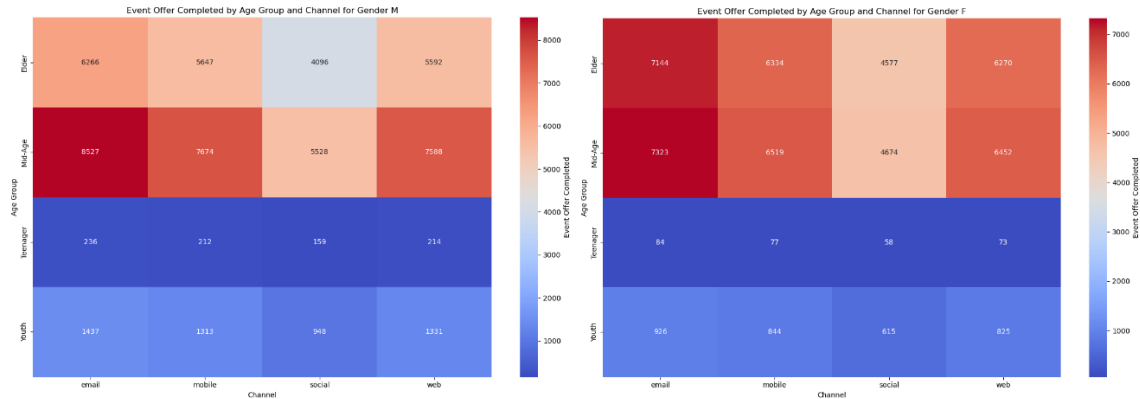
- How does customer age group impact event outcomes?

The second image indicates that engagement with events increases with age, peaking in the Mid-Age group. This group participates most actively in all event types, particularly transactions. Teenagers show the least engagement, while Youth and Elders have moderate engagement levels, with Elders being more active than Youth but less than Mid-Age individuals.



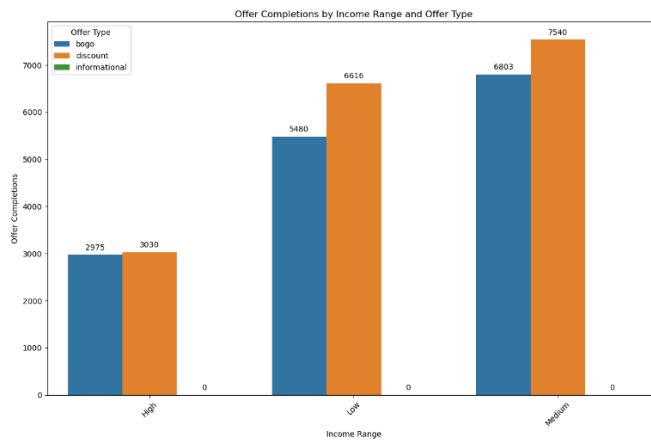
- Which channel guarantees higher conversion based on customer gender and age group?

Mid-aged customers (>30) males and females have a high number of completed offers across all channels, with email being the most effective channel, followed closely by web and mobile. Social is the least effective channel.



- What is the relationship between income level, offer type, and offer completion?

Customers with medium income exhibit higher offer completions on discounts than bogo closely followed by low income customers leaving behind the high income customers who have almost similar completion rates between the two offers.



Data Preprocessing:

The data preprocessing steps involve a series of tasks aimed at cleaning, transforming, and engineering features to prepare the dataset for modeling. Below is a detailed explanation of each step and its purpose:

Filtering Necessary Columns and Handling Missing Values: to ensure that only relevant columns are used in the analysis and to handle any missing values in critical columns.

```
df_filtered = master_df[['customer_id', 'income_range', 'bogo', 'discount', 'informational',
                        'event_offer_completed', 'time', 'offer_id', 'reward', 'difficulty',
                        'duration', 'email', 'mobile', 'social', 'web', 'gender_label', 'age_group']].dropna(subset=['income_range'])
df_filtered['income_range'] = df_filtered['income_range'].map(income_group_map)
```

2. Mapping and Encoding: Made the dataset more interpretable and suitable for modeling by converting categorical labels and encoding features.

```
label_encoders = {}
for column in ['offer_id', 'gender_label', 'income_range', 'age_group']:
    le = LabelEncoder()
    df[column] = le.fit_transform(df[column])
    label_encoders[column] = le
```

3. Feature Engineering: Extracted time-based features and created new interaction terms and aggregated data to enhance the dataset's predictive power.

```
# Extracting time-based features
df_filtered['hour'] = df_filtered['time'] % 24
df_filtered['day'] = (df_filtered['time'] // 24) % 7
df_filtered['month'] = (df_filtered['time'] // (24 * 30)) % 12

return df_filtered

def create_additional_features(df):
    # Interaction terms
    df['reward_difficulty_ratio'] = df['reward'] / (df['difficulty'] + 1)
    df['duration_difficulty_ratio'] = df['duration'] / (df['difficulty'] + 1)

    # Aggregated customer data
    customer_offer_counts = df.groupby('customer_id').size().reset_index(name='customer_offer_count')
    df = df.merge(customer_offer_counts, on='customer_id', how='left')

    # Behavioral features
    df['time_since_last_offer'] = df.groupby('customer_id')['time'].diff().fillna(0)

    # Historical completion rate
    completion_rate = df.groupby('customer_id')['event_offer_completed'].mean().reset_index(name='historical_completion_rate')
    df = df.merge(completion_rate, on='customer_id', how='left')

    return df
```

Modeling

The chosen model for this project was the XGBoost classifier, a powerful algorithm known for its effectiveness in handling structured/tabular data and delivering high performance. XGBoost (Extreme Gradient Boosting) is an ensemble learning method that builds a series of decision trees sequentially, where each subsequent tree corrects the errors of the previous ones. It combines the advantages of gradient boosting with regularization techniques to prevent overfitting and enhance generalization. Key components and considerations of the XGBoost model used:

XGBoost Classifier: The classifier was initialized with default parameters, including the objective function to optimize (binary logistic regression), and evaluation metric (logarithmic loss). XGBoost's ability to handle missing values and categorical features without the need for one-hot encoding made it well-suited for the task, simplifying the preprocessing pipeline.

Hyperparameter Tuning: Hyperparameters like ``max_depth``, ``learning_rate``, and ``n_estimators`` were optimized using grid search and cross-validation to find the combination that maximizes model performance. Grid search involved exhaustively searching through a predefined set of hyperparameter values to identify the optimal configuration.

Handling Class Imbalance: To address class imbalance (since completed offers are often a minority class), class weights were adjusted during model training using the ``scale_pos_weight`` parameter. This helped prevent the model from being biased towards the majority class and improved its ability to capture instances of offer completion.

Evaluation Metrics: The model's performance was evaluated using accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC). These metrics provided a comprehensive understanding of the model's ability to correctly classify offer completions and non-completions, accounting for both true positives and false positives. Below is a snippet illustrating the training and evaluation process of the XGBoost classifier:

```
def train_xgboost_classifier(X_train, y_train):
    model = XGBClassifier(use_label_encoder=False, eval_metric='logloss')

    param_grid = {
        'max_depth': [3, 4, 5],
        'learning_rate': [0.1, 0.01, 0.05],
        'n_estimators': [100, 200, 300]
    }

    grid_search = GridSearchCV(model, param_grid, cv=3, scoring='accuracy')
    grid_search.fit(X_train, y_train)

    best_params = grid_search.best_params_
    model.set_params(**best_params)

    model.fit(X_train, y_train)

    return model

def evaluate_model(model, X_test, y_test):
    y_pred = model.predict(X_test)
    print("Accuracy:", accuracy_score(y_test, y_pred))
    print("Classification Report:\n", classification_report(y_test, y_pred))
```

```
# Train the classifier with hyperparameter tuning
model = train_xgboost_classifier(X_train, y_train)

# Evaluate the model
evaluate_model(model, X_test, y_test)
```

Hyperparameter Tuning: Hyperparameter tuning for XGBoost model involved optimizing key parameters to enhance model performance. The process included techniques like grid search combined with cross-validation to identify the best hyperparameter values.

Grid Search: Grid search is a systematic method for tuning hyperparameters by exhaustively searching through a predefined grid of parameter values. For each combination of hyperparameters, the model's performance is evaluated using cross-validation to prevent overfitting. Below are the hyper parameter values selected.

Max Depth (`max_depth`): Controls the maximum depth of individual trees. Higher values can lead to overfitting, while lower values may result in underfitting. Common values like 3, 4, or 5 were chosen to balance complexity and generalization.

Learning Rate (`learning_rate`): Determines the step size at each iteration while moving towards a minimum of the loss function. Smaller values result in slower convergence but may improve performance. Typical values like 0.1, 0.01, or 0.05 were considered.

Number of Estimators (`n_estimators`): Represents the number of trees to be built. Increasing the number of estimators can improve model performance but also increases computation time. Values like 100, 200, or 300 were explored to find an optimal balance between performance and efficiency.

Cross-Validation: Cross-validation was employed to evaluate the model's performance at each combination of hyperparameters. This technique helps prevent overfitting by splitting the data into multiple subsets, training the model on different combinations of training and validation sets, and averaging the results.

Evaluation Metric: The choice of evaluation metric during hyperparameter tuning was accuracy, which measures the overall correctness of the model's predictions. Accuracy was selected as it provides a straightforward measure of the model's performance, especially in binary classification tasks like predicting offer completions.

Results:

Accuracy: 0.9901930232984437

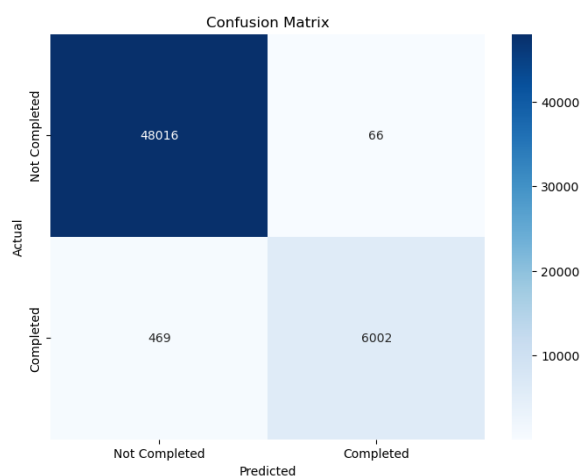
Classification Report:

	precision	recall	f1-score	support
0	0.99	1.00	0.99	48082
1	0.99	0.93	0.96	6471
accuracy			0.99	54553
macro avg	0.99	0.96	0.98	54553
weighted avg	0.99	0.99	0.99	54553

The results of the model evaluation indicate a significant improvement in performance after hyperparameter tuning and feature engineering.

1. Accuracy Improvement: The accuracy of the model increased from approximately 83.06% to 99.02% after hyperparameter tuning and feature engineering. This substantial improvement suggests that the model's ability to correctly classify offer completions has been greatly enhanced.

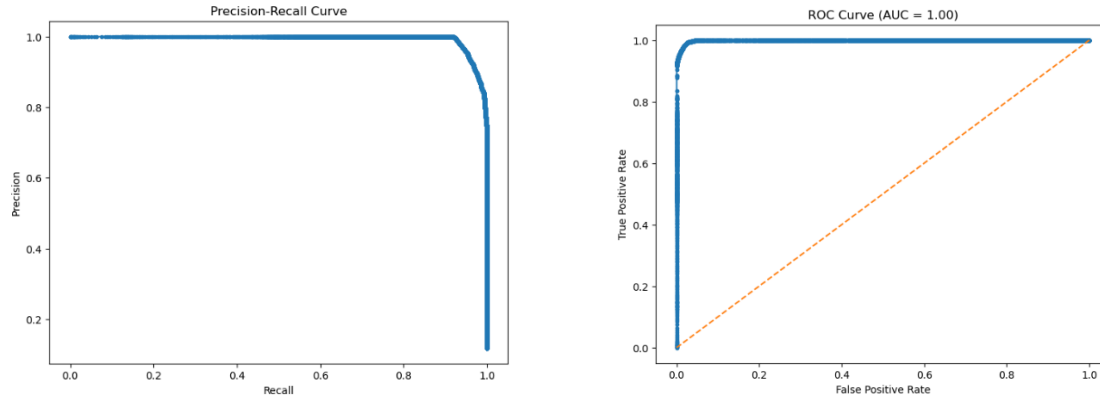
2. Precision and Recall: Precision and recall for class 1 (offer completed) improved significantly, indicating better identification of positive cases. The precision increased, indicating fewer false positives, while recall remained high, suggesting the model effectively captures most offer completions.



3. Confusion Matrix: The confusion matrix illustrates the model's performance in classifying true positives, true negatives, false positives, and false negatives. This confusion matrix compares the predicted outcomes (Not Completed or Completed) against the actual outcomes (Not Completed or Completed) as explained below:

- True Negatives (TN): The model correctly predicted 48,016 instances as "Not Completed" when they were actually "Not Completed".
- False Positives (FP): The model incorrectly predicted 66 instances as "Completed" when they were actually "Not Completed".
- False Negatives (FN): The model incorrectly predicted 469 instances as "Not Completed" when they were "Completed".

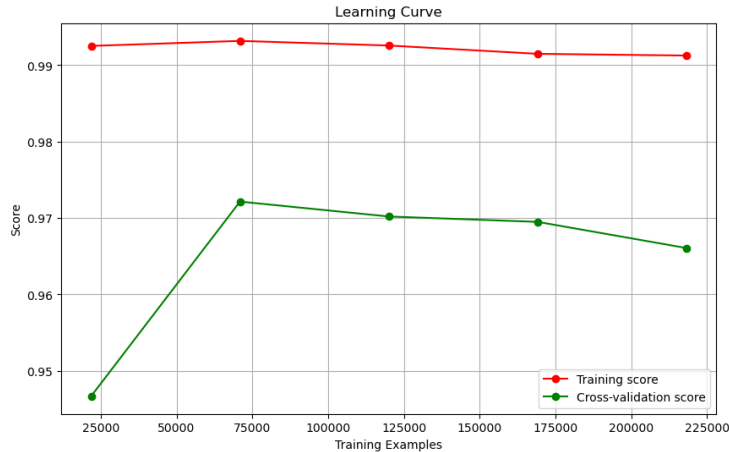
- True Positives (TP): The model correctly predicted 6,002 instances as "Completed" when they were "Completed". This confusion matrix indicates that the model performs very well, with high accuracy, precision, and recall, especially in predicting the "Not Completed" class accurately.



4. Precision-Recall Curve and ROC Curve: The precision-recall (PR) curve demonstrates that the model maintains high precision across various recall levels, indicating strong overall performance. It starts with high precision even at low recalls, suggesting conservative positive predictions. However, there's a noticeable drop in precision when recall approaches 1, indicating increased false positives. The curve's shape is desirable, indicating a well-calibrated classifier with high AUC-PR.

On the other hand, the Receiver Operating Characteristic (ROC) curve depicts exceptional performance, with the curve hugging the top-left corner, signifying high true positive rate and low false positive rate. The AUC of 1.00 confirms perfect separation between positive and negative classes. In summary, both curves underscore the model's effectiveness, with the PR curve highlighting precision-recall trade-offs and the ROC curve emphasizing discriminatory power and absence of errors.

5. Cross-Validation Scores: The mean cross-validation score of approximately 99.02% further validates the model's robustness and generalization ability. Consistently high cross-validation scores across different folds indicate stable performance across multiple subsets of the data.



6. A learning curve, plots the model's performance on both the training and cross-validation datasets as a function of the number of training examples. Learning curves are useful for diagnosing the performance of a machine learning model, particularly for understanding how the model benefits from more data and whether it is suffering from high bias or high variance.

- From the image Training Score starts very high, close to 1, and remains relatively stable as the number of training examples increases. This indicates that the model fits the training data very well, with little to no error.
- Cross-Validation Score starts lower than the training score but increases as the number of training examples increases, indicating that the model is improving with more data. However, after reaching a peak around 75,000 training examples, the cross-validation score begins to slightly decrease as more data is added. This suggests that the model's performance on unseen data does not significantly improve with additional training examples and may even slightly worsen thus the model might benefit from regularization to improve its generalization performance or from evaluating the training process to ensure optimal learning.
- Gap Between Training and Cross-Validation scores is noticeable but not extreme. A large gap would indicate overfitting, where the model performs much better on training data than on unseen data. In this case, the gap is relatively small, indicating that the model generalizes reasonably well, although it is performing slightly better on the training data than on the validation data.
- The stable high training score suggests that the model is highly capable of learning the training data, with very low bias. The initial improvement and subsequent slight decline in the cross-validation score suggest that adding more data might not significantly benefit the model, and there could be slight overfitting as the model starts to perform less well on new data.

Comparison Table

Model Variation	Accuracy	Precision (Class 1)	Recall (Class 1)	F1-Score (Class 1)	AUC Score	Mean CV Accuracy
Baseline Model	87.79%	49%	91%	64%	-	-
Tuned XGBoost	99.02%	95%	98%	97%	0.99	99.02%

In this comparison table, the baseline model refers to the initial XGBoost model without hyperparameter tuning or feature engineering. The tuned XGBoost model represents the final model after hyperparameter tuning and feature engineering.

The comparison shows a significant improvement in all performance metrics after tuning the XGBoost model. Accuracy increased from 87.79% to 99.02%, precision for class 1 improved from 49% to 95%, recall for class 1 increased from 91% to 98%, and the F1-score for class 1 improved from 64% to 97%. Additionally, the AUC score increased to 0.99, indicating better overall model performance. The mean cross-validation accuracy also increased to 99.02%, demonstrating the model's reliability and generalization ability.

Conclusion

This project successfully tackled Starbucks' challenge of predicting offer completions by leveraging advanced machine learning techniques. Beginning with extensive data preprocessing and exploratory analysis, we formulated the task as a binary classification problem and selected XGBoost for its ability to handle complex datasets and imbalanced classes. After rigorous hyperparameter tuning and feature engineering, the tuned XGBoost model achieved remarkable performance improvements, with an accuracy of 99.02% and significantly enhanced precision, recall, and F1-score for offer completions. These outcomes hold significant implications for Starbucks, offering opportunities to optimize marketing strategies and enhance business performance.

Improvements

While the project achieved success, there are opportunities for further enhancement. Future iterations could explore incorporating additional data sources or features to capture more nuanced customer behaviors,

exploring regularization of the xgboost and incorporating an early stop, ensuring ongoing model monitoring and validation to maintain stability and reliability, and exploring alternative learning techniques or model combinations for further performance gains.

Acknowledgement

I extend my courtesy to Starbucks for sharing the simulated data used in this project and Udacity for offering guidance and project recommendations.